

# FMAN-45 Machine Learning, Spring 2019

## Assignment 3: Convolutional Neural Networks

*Solve the problems and write down the solutions. Download the code we provide and do the additional, required programming. Write a detailed report. All solutions, plots and figures should be in one pdf. It should be possible to understand all material presented in the report without running any code. Submit your solutions and code using your individual Moodle account as a single archive at <http://moodle.maths.lth.se/course/> by the deadline (note that there will be no extensions). Do not include the `data` directory in the archive.*

### 1 Common Layers and Backpropagation

In this section you will implement forward and backward steps for a few common neural network layers.

When you train a neural net you have to evaluate the gradient  $\frac{\partial L}{\partial \theta}$  with respect to all the parameters in the network. The algorithm to do this is called backpropagation and it is illustrated in figure 1.

#### 1.1 Fully Connected Layer

Suppose that we have a vector  $\mathbf{x} \in \mathbb{R}^n$  and we define a mapping to  $\mathbf{y} \in \mathbb{R}^m$  by

$$y_i = \sum_{j=1}^n A_{ij}x_j + b_i. \quad (1)$$

Using matrix notation we have  $\mathbf{y} = \mathbf{Ax} + \mathbf{b}$ . Note that  $\mathbf{A}$  and  $\mathbf{b}$  are parameters that are trainable. Comparing with figure 1,  $\mathbf{x}$  might correspond to  $\mathbf{x}_1$ ,  $\mathbf{y}$  to  $\mathbf{x}_2$  and the parameter  $\theta_2 = \{\mathbf{A}, \mathbf{b}\}$ . Thus several such mappings, or *layers*, are concatenated to form a full network such as the one shown in figure 1. To derive the equations for backpropagation we want to express  $\frac{\partial L}{\partial \mathbf{x}}$  as an expression containing  $\frac{\partial L}{\partial \mathbf{y}}$ . Using the chain rule we get

$$\frac{\partial L}{\partial x_i} = \sum_{j=1}^m \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial x_i}. \quad (2)$$

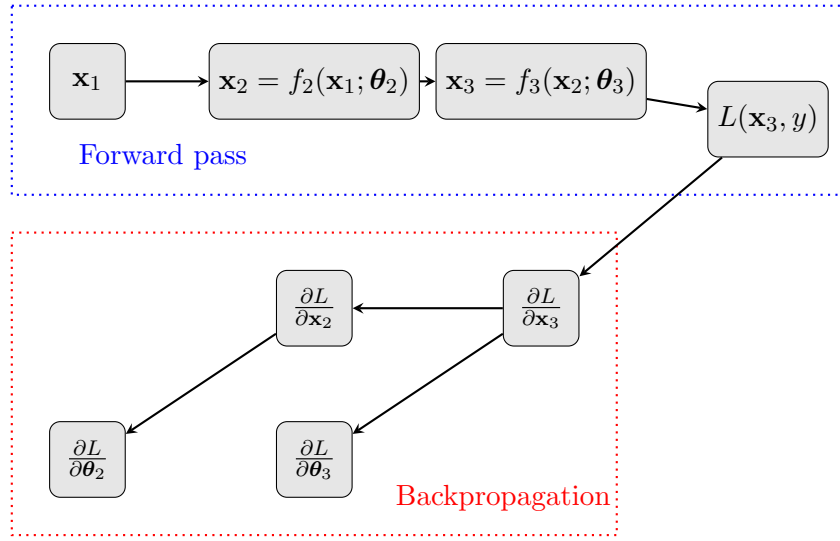


Figure 1: Illustration of the computational graph for a simple network with two layers and a loss. Given an input  $\mathbf{x}_1$  we first do the forward pass and compute the layer activations  $\mathbf{x}_2$ ,  $\mathbf{x}_3$  and the loss  $L$ . The loss depends on the ground truth  $y$ . When we go in the other direction the goal is to obtain the gradients with respect to the parameters  $\boldsymbol{\theta}_2$  and  $\boldsymbol{\theta}_3$ . First you compute the gradients with respect to the layer activations  $\frac{\partial L}{\partial \mathbf{x}_3}$  and  $\frac{\partial L}{\partial \mathbf{x}_2}$ . Finally the gradients with respect to the parameters  $\frac{\partial L}{\partial \boldsymbol{\theta}_3}$  and  $\frac{\partial L}{\partial \boldsymbol{\theta}_2}$  can be obtained from  $\frac{\partial L}{\partial \mathbf{x}_3}$  and  $\frac{\partial L}{\partial \mathbf{x}_2}$  respectively.

We also need to compute the gradient with respect to the parameters  $\mathbf{A}$  and  $\mathbf{b}$ . To do this, again use the chain rule,

$$\frac{\partial L}{\partial A_{ij}} = \sum_{k=1}^m \frac{\partial L}{\partial y_k} \frac{\partial y_k}{\partial A_{ij}} \quad (3)$$

$$\frac{\partial L}{\partial b_i} = \sum_{j=1}^m \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial b_i} \quad (4)$$

**Exercise 1 (10 points):** Give the simplified expressions for  $\frac{\partial L}{\partial \mathbf{x}}$ ,  $\frac{\partial L}{\partial \mathbf{A}}$  and  $\frac{\partial L}{\partial \mathbf{b}}$ . Give a full solution. The answers should all be given as matrix expressions without any explicit sums.

When you have code that uses gradients, it is very important to check that the gradients are correct. A bad way to conclude that the gradient is correct is to manually look at the code and convince yourself that it is correct or just to see if the function seems to decrease when you run optimization; it might still be a descent direction even if it is not the gradient. A much better way is to check finite differences using the formula

$$\frac{\partial f}{\partial x_i} \approx \frac{f(\mathbf{x} + \epsilon \mathbf{e}_i) - f(\mathbf{x} - \epsilon \mathbf{e}_i)}{2\epsilon}, \quad (5)$$

where  $\mathbf{e}_i$  is a vector that is all zero except for position  $i$  where it is 1, and  $\epsilon$  is a small number. In the code there is a file `tests/test_fully_connected.m` where you can test your implementation.

When you are training a neural network it is common to evaluate the network and compute gradients not with respect to just one element but  $N$  elements in a batch. We use superscripts to denote the elements in the batch. For the fully connected layer above, we have multiple inputs  $\mathbf{x}^{(1)}$ ,  $\mathbf{x}^{(2)}$ ,  $\dots$ ,  $\mathbf{x}^{(N)}$  and we wish to compute  $\mathbf{y}^{(1)} = \mathbf{A}\mathbf{x}^{(1)} + \mathbf{b}$ ,  $\mathbf{y}^{(2)} = \mathbf{A}\mathbf{x}^{(2)} + \mathbf{b}$ ,  $\dots$ ,  $\mathbf{y}^{(N)} = \mathbf{A}\mathbf{x}^{(N)} + \mathbf{b}$ . In the code for the forward pass you can see that the input array first is reshaped to a matrix  $\mathbf{X}$  where each column contains all values for a single batch, that is,

$$\mathbf{X} = (\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots \quad \mathbf{x}^{(N)}). \quad (6)$$

For instance, if  $\mathbf{x}^{(1)}$  an image of size  $5 \times 5 \times 3$  it is reshaped to a long vector with length 75. We wish to compute

$$\mathbf{Y} = (\mathbf{y}^{(1)} \quad \mathbf{y}^{(2)} \quad \dots \quad \mathbf{y}^{(N)}) = (\mathbf{A}\mathbf{x}^{(1)} + \mathbf{b} \quad \mathbf{A}\mathbf{x}^{(2)} + \mathbf{b} \quad \dots \quad \mathbf{A}\mathbf{x}^{(N)} + \mathbf{b}). \quad (7)$$

When we are backpropagating to  $\mathbf{X}$  we have to compute

$$\frac{\partial L}{\partial \mathbf{X}} = \left( \frac{\partial L}{\partial \mathbf{x}^{(1)}} \quad \frac{\partial L}{\partial \mathbf{x}^{(2)}} \quad \dots \quad \frac{\partial L}{\partial \mathbf{x}^{(N)}} \right) \quad (8)$$

using

$$\frac{\partial L}{\partial \mathbf{Y}} = \left( \frac{\partial L}{\partial \mathbf{y}^{(1)}} \quad \frac{\partial L}{\partial \mathbf{y}^{(2)}} \cdots \frac{\partial L}{\partial \mathbf{y}^{(N)}} \right). \quad (9)$$

Use your expression obtained in exercise 1 and simplify to matrix operations. For the parameters, we have that both  $\mathbf{A}$  and  $\mathbf{b}$  influence all elements  $\mathbf{y}^{(i)}$ , so for the parameters we compute  $\frac{\partial L}{\partial \mathbf{A}}$  using the chain rule with respect to each element in each  $\mathbf{y}^{(i)}$  and just sum them up. More formally,

$$\frac{\partial L}{\partial A_{ij}} = \sum_{l=1}^N \sum_{k=1}^m \frac{\partial L}{\partial y_k^{(l)}} \frac{\partial y_k^{(l)}}{\partial A_{ij}} \quad (10)$$

$$\frac{\partial L}{\partial b_i} = \sum_{l=1}^N \sum_{j=1}^m \frac{\partial L}{\partial y_j^{(l)}} \frac{\partial y_j^{(l)}}{\partial b_i}. \quad (11)$$

Note that the inner sum is the same as you have computed in the previous exercise, so you just sum the expression you obtained in the previous exercise over all elements in the batch. You can implement the forward and backward passes with for-loops, but it is also possible to use matrix operations to vectorize the code and for full credit you should vectorize it. Potentially useful functions in Matlab include `bsxfun`, `sub2ind`, `ind2sub`, `reshape` and `repmat`.

**Exercise 2 (10 points):** Write a full solution of how  $\mathbf{Y}$ ,  $\frac{\partial L}{\partial \mathbf{X}}$ ,  $\frac{\partial L}{\partial \mathbf{A}}$  and  $\frac{\partial L}{\partial \mathbf{b}}$  should be vectorized (do *not* just copy your code for vectorization. Give a mathematical proof of why it works). In the code there are two files `layers/fully_connected_forward.m` and `layers/fully_connected_backward.m`. Implement these functions and check that your implementation is correct by running `tests/test_fully_connected.m`. For full credit your code should be vectorized over the batch. Include the relevant code in the report.

## 1.2 Relu

The most commonly used function as a nonlinearity is the rectified linear unit (relu). It is defined as

$$y_i = \max(x_i, 0) \quad (12)$$

**Exercise 3 (10 points):** Derive the backpropagation expression for  $\frac{\partial L}{\partial x_i}$ . Give a full solution. Implement the layer in `layers/relu_forward.m` and `layers/relu_backward.m`. Test it with `tests/test_relu.m`. For full credit you must use built in indexing and not use any for-loops. Include the relevant code in the report.

### 1.3 Softmax Loss

Suppose we have a vector  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$ . The goal is to classify the input as one out of  $n$  classes and  $x_i$  is a score for class  $i$  and a larger  $x_i$  is a higher score. By using the softmax function we can interpret the scores  $x_i$  as probabilities. Let

$$y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \quad (13)$$

be the probability for class  $i$ . Now suppose that the ground truth class is  $c$ . Since now  $y_i$  are probabilities we can define the loss  $L$  to be minimized as the negative log likelihood,

$$L(\mathbf{x}, c) = -\log(y_c) = -\log\left(\frac{e^{x_c}}{\sum_{j=1}^n e^{x_j}}\right) = -x_c + \log\left(\sum_{j=1}^n e^{x_j}\right). \quad (14)$$

**Exercise 4 (10 points):** Compute the expression for  $\frac{\partial L}{\partial x_i}$ . Write a full solution. Note that this is the final layer, so there are no gradients to backpropagate. Implement the layer in `layers/softmaxloss_forward.m` and `layers/softmaxloss_backward.m`. Test it with `tests/test_softmaxloss.m`. Make sure that `tests/test_gradient_whole_net.m` runs correctly when you have implemented all layers. For full credit you should use built in functions for summation and indexing and not use any for-loops. Include the relevant code in the report.

## 2 Training a Neural Network

The function we are trying to minimize when we are training a neural net is

$$L(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta}) \quad (15)$$

where  $\boldsymbol{\theta}$  are all the parameters of the network,  $\mathbf{x}^{(i)}$  is the input and  $y^{(i)}$  the corresponding ground truth and  $L(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta})$  is the loss for a single example, given for instance by a neural net with a softmax loss in the last layer, as in figure 1. In practice,  $N$  is very large and we never evaluate the gradient over the entire sum. Instead we evaluate it over a batch with  $n \ll N$  elements because the network can be trained much faster if you use batches and update the parameters in every step.

If you are using gradient descent to train the network, the way the parameters are updated is

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \alpha \frac{\partial L}{\partial \boldsymbol{\theta}} \quad (16)$$

where  $\alpha$  is a hyperparameter called the learning rate. Since we only evaluate the gradient over a few examples, the estimated gradient in (16) might be very noisy. The idea behind gradient descent with momentum is to average the gradient estimations over time and use the smoothed gradient to update the parameters. The update in (16) is modified as

$$\mathbf{m}_n = \mu \mathbf{m}_{n-1} + (1 - \mu) \frac{\partial L}{\partial \boldsymbol{\theta}} \quad (17)$$

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \alpha \mathbf{m}_n \quad (18)$$

where  $\mathbf{m}_n$  is a moving average of the gradient estimations and  $\mu$  is a hyperparameter in the range  $0 < \mu < 1$  controlling the smoothness.

**Exercise 5 (10 points):** Implement gradient descent with momentum in `training.m`. Remember to include weight decay. Include the relevant code in the report.

### 3 Classifying Handwritten Digits

**You must have solved all previous problems before you can start working on this and the next problem.**

In `mnist_starter.m` there is a simple baseline for the MNIST dataset. Read the comments in the code carefully. It reaches about 98 % accuracy on the test set (this of course varies a bit from time to time). Validate this to make sure that the code you have written so far is correct.

**Exercise 6 (25 points):** Plot the filters the first convolutional layer learns. Plot a few images that are misclassified. Plot the confusion matrix for the predictions on the test set and compute the precision and the recall for all digits. Write down the number of parameters for all layers in the network. Write comments about all plots and figures.

### 4 Classifying Tiny Images

In `cifar10_starter.m` is a simple network that can be trained on the cifar10 dataset. This baseline give an accuracy of about 48 % after training for 5000 iterations. Note that it is much harder to get good classification results on this dataset than MNIST, mainly due to significant intraclass variation.

**Exercise 7 (25 points):** Do whatever you want to improve the accuracy of the baseline network. Write what you have tried in the report, even experiments that did not work out. For your final model, compute and report all plots and numbers that were required in the previous exercise and comment on it.