

# Simple Convolutional Neural Network From Scratch

Chris Yogodzinski

## Notebook Overview

This notebook is a follow up to `src/python/scratch_two_layer.ipynb`.

The previous notebook I made a simple neural network from scratch to predict handwritten digits in the mnist dataset. The neural network consisted of one fully connected hidden layer and an output layer.

In this notebook I will be adding a convolutional layer and pooling layer to this simple neural network. The convolutions in this neural network function to key features from the dataset prior to the predictive fully connected layer.

## Variables

$I$  matrix: a mnist image of dim = ( $m$  by  $n$ )

$K$  matrix: kernel of dim = ( $m_k$  by  $n_k$ )

$s$  integer: stride, a constant set by user

$\phi_a$  function: activation function of convolution

$F$  matrix: a feature map of dim = ( $m - m_k + 1$  by  $n - n_k + 1$ ) or ( $i$  by  $j$ )

$\phi_p$  function: pooling function applied to feature map

$P$  matrix: pooled feature map, condensed image of dim = ( $\frac{m-m_k}{s}$  by  $\frac{n-n_k}{s}$ )

## Convolutional Layer

The convolutional layer in a cnn is named after a method known as kernel convolution.

Kernel convolution works by passing a filter or kernel iteratively over the input matrix and transform the values of the input based on the values of the kernel.

In this notebook our kernel transformation will be  $\sum_i \sum_j I[m-i, n-j] \times K$ .

In words this is element-wise multiplication of the kernel-sized subset of  $I$  and  $K$  followed by the sum of the product.

Kernel Convolution simplifies our image which reduces the number of parameters our model needs to learn in the fully connected layers.

However the kernel itself is also a set of parameters that needs to be learned so that we're not losing important information when we simplify the images.

## Importing the data

The following code chunks read in the binary formatted mnist images

```
setwd("~/Projects/nn_playground")

read_mnist <- function(label_path, im_path) {
  # read in labels
```

```

f <- file(label_path, "rb")
meta <- readBin(f, n = 2, "integer", endian = "big")
labels <- as.integer(readBin(f, n = meta[2], "raw", endian = "big"))
close(f)

# read in imgs
f <- file(im_path, "rb")
meta <- readBin(f, n = 4, "integer", endian = "big")
byte_count <- meta[2] * meta[3] * meta[4]
imgs <- readBin(f, n = byte_count, "raw", endian = "big")
close(f)
imgs <- array(as.integer(imgs), dim = c(meta[3], meta[4], meta[2]))
dat <- list(labels = labels, imgs = imgs)
return(dat)
}

dat_path = "data/mnist"
train_dat <- read_mnist(
  paste(dat_path, "train-labels.idx1-ubyte", sep = "/"),
  paste(dat_path, "train-images.idx3-ubyte", sep = "/")
)

test_dat <- read_mnist(
  paste(dat_path, "t10k-labels.idx1-ubyte", sep = "/"),
  paste(dat_path, "t10k-images.idx3-ubyte", sep = "/")
)

```

## Image Preprocessing

### Flipping The Images

The function above reads the images in upside down for some reason.

In the code chunk below I flip the images right side up and visualize them. `###` Adding padding to images  
 In this step I'm adding a border of empty pixels around the edge of the images.

The purpose of this step is to allow our model to get a better “view” of the features around the edges of the image. It is likely this is unnecessary for such simple images like the mnist dataset. In this step I am also scaling the pixel values to be between 0 and 1.

```

preprocess_img <- function(img) {

  flip_img <- function(img) {
    flipped <- img[, ncol(img):1]
    return(flipped)
  }

  zero_pad_img <- function(img) {
    top_bottom <- rep_len(0, ncol(img))
    img <- rbind(top_bottom, img, top_bottom)
    left_right <- rep_len(0, nrow(img))
    img <- unname(cbind(left_right, img, left_right))
    return(img)
  }

  preprocessed_img <- zero_pad_img(flip_img(img)) / 255
}

```

```

    return(preprocessed_img)
}

preprocess_all <- function(img_array) {
  dims <- dim(img_array) + c(2, 2, 0)
  preprocessed <- array(apply(img_array, 3, preprocess_img), dim = dims)
  return(preprocessed)
}

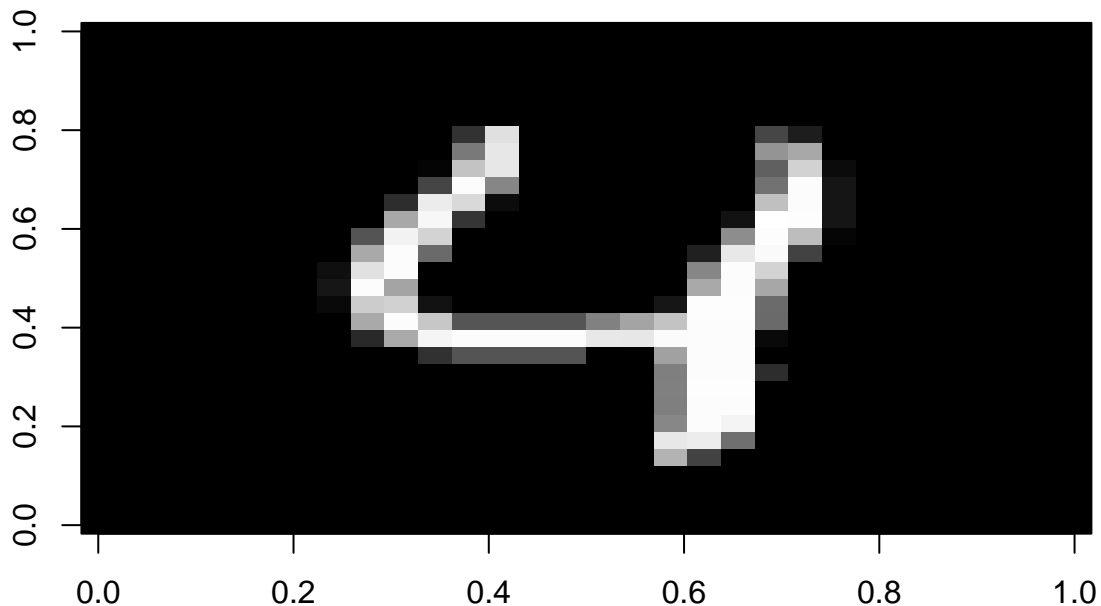
# apply preprocessing
train_dat[["imgs"]] <- preprocess_all(train_dat[["imgs"]])
test_dat[["imgs"]] <- preprocess_all(test_dat[["imgs"]])

im_num <- 5
# test image
img <- test_dat[["imgs"]][ , , im_num]
cat(test_dat[["labels"]][im_num], sep = "\n")

```

```
## 4
```

```
image(img, col = gray((0:255) / 255))
```



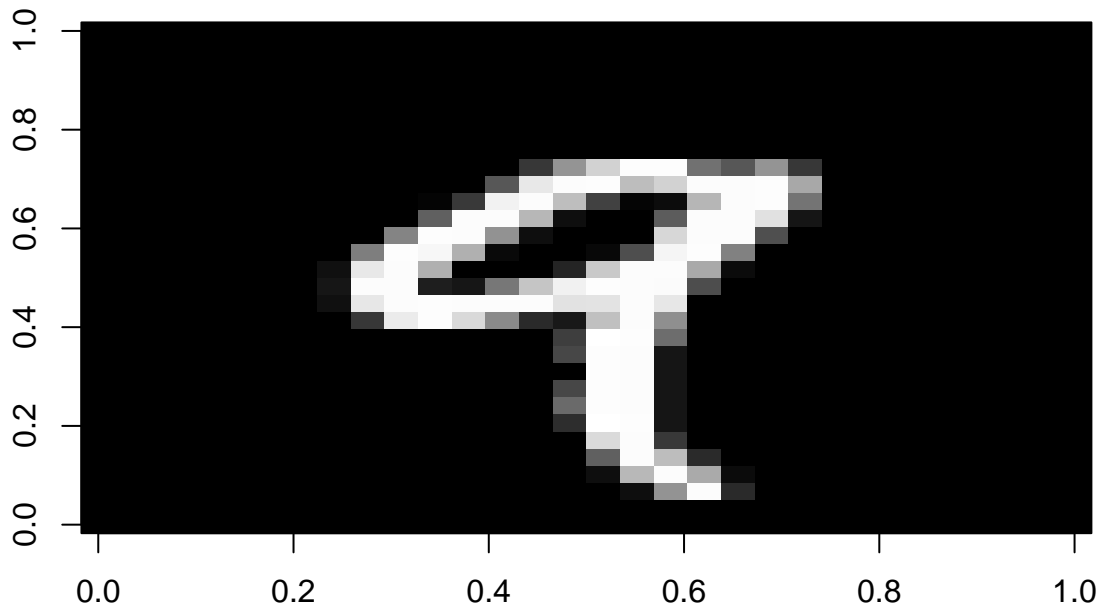
```

# train image
img <- train_dat[["imgs"]][ , , im_num]
cat(train_dat[["labels"]][im_num], sep = "\n")

```

```
## 9
```

```
image(img, col = gray((0:255) / 255))
```



## Forward Propagation

In the following cells i've written some functions needed for forward propagation

### Kernel Convolution

In the cell below I've written some functions to apply kernel convolution to an array of our images. You can see from the images printed at below the cell how some features of the digits have been filtered out after kernel convolution.

Initially I implemented kernel convolution in the spatial domain. Using this procedure kernel convolution of the complete training dataset would take roughly 40 mins.

I could implement some parallel processing or train on a gpu to increase performance.

Instead I decided to implement kernel convolution in the spatial domain using fast fourier transform. This drastically improves performance.

```
kernel_conv_slow <- function(img, kernel, kernel_dim, feat_map_idx, dim_out) {

  # wrapper function to filter signal panel of image
  kernel_filt <- function(x) {
    row_idx <- feat_map_idx[x, "Var1"]
    col_idx <- feat_map_idx[x, "Var2"]
    img_panel <- img[row_idx:(row_idx+kernel_dim[1] - 1), col_idx:(col_idx+kernel_dim[2] - 1)]
    filtered_panel <- sum(img_panel * kernel)
    return(filtered_panel)
  }

  # Apply complete convolution
  feat_map <- matrix(
    sapply(rownames(feat_map_idx), FUN = kernel_filt),
    nrow = dim_out,
    ncol = dim_out
  )
  return(feat_map)
}
```

```

# pads kernel with zeros to be same dimensions as image, needed for fft convolution
pad_kernel <- function(kern, img_dim) {
  padding_len <- img_dim[1] - dim(kern)[1]
  right <- matrix(0, nrow = dim(kern)[1], ncol = padding_len)
  bottom <- matrix(0, nrow = padding_len, ncol = img_dim[2])
  padded <- rbind(cbind(kern, right), bottom)
  return(padded)
}

kernel_conv_fft <- function(img, img_dims, padded_ft, dim_out) {
  convolved <- fft(fft(img) * padded_ft, inverse = TRUE)
  return(convolved)
}

apply_kernel_conv <- function(imgs, kern, kernel_dim, stride, method = "fft") {

  data_dim <- dim(imgs)
  dim_out <- ((data_dim[1] - kernel_dim[1]) / stride) + 1

  if (method == "fft") {
    # zero pad kernel to match size of images
    padded <- pad_kernel(kern[kernel_dim[1]:1, kernel_dim[2]:1], data_dim[1:2])
    padded_ft <- fft(padded)
    dim_out <- data_dim
    conv_func <- function(x) {
      kernel_conv_fft(x, data_dim[1:2], padded_ft, dim_out)
    }
  } else {
    # generate two column matrix of all indeces needed for full convolution
    conv_idx <- seq(1, data_dim[2] - kernel_dim[2] + 1, by = stride)
    feat_map_idx <- expand.grid(rep(list(conv_idx), 2))
    dim_out <- c(dim_out, dim_out, data_dim[3])
    conv_func <- function(x) {
      kernel_conv_slow(x, kern, kernel_dim, feat_map_idx, dim_out)
    }
  }
  feat_maps <- apply(imgs, MARGIN = 3, conv_func)
  feat_maps <- array(as.numeric(feat_maps), dim = dim_out)
  return(feat_maps)
}

```

Below I have benchmarked the two convolution implementations. I also show the differences in the images after convolution.

```

stride <- 1
kernel_dim <- c(3, 3)
kern <- matrix(
  rnorm(prod(kernel_dim), 0, 1.5),
  nrow = kernel_dim[1],
  ncol = kernel_dim[2]
)

```

```
ti <- Sys.time()
feature_maps_slow <- apply_kernel_conv(train_dat[["imgs"]][, , 1:100], kern, kernel_dim, stride, method = "slow")
tf <- Sys.time()
print(tf - ti)
```

```
## Time difference of 12.34986 secs
```

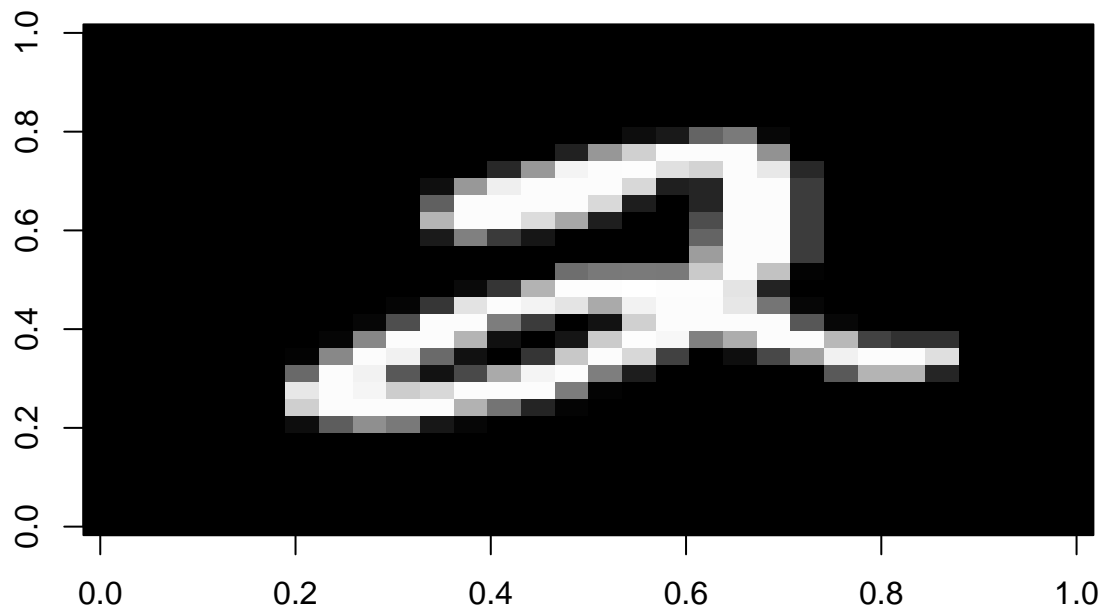
```
ti <- Sys.time()
feature_maps_fft <- apply_kernel_conv(train_dat[["imgs"]][, , 1:100], kern, kernel_dim, stride, method = "fft")
tf <- Sys.time()
print(tf - ti)
```

```
## Warning in array(as.numeric(feats_maps), dim = dim_out): imaginary parts
## discarded in coercion
```

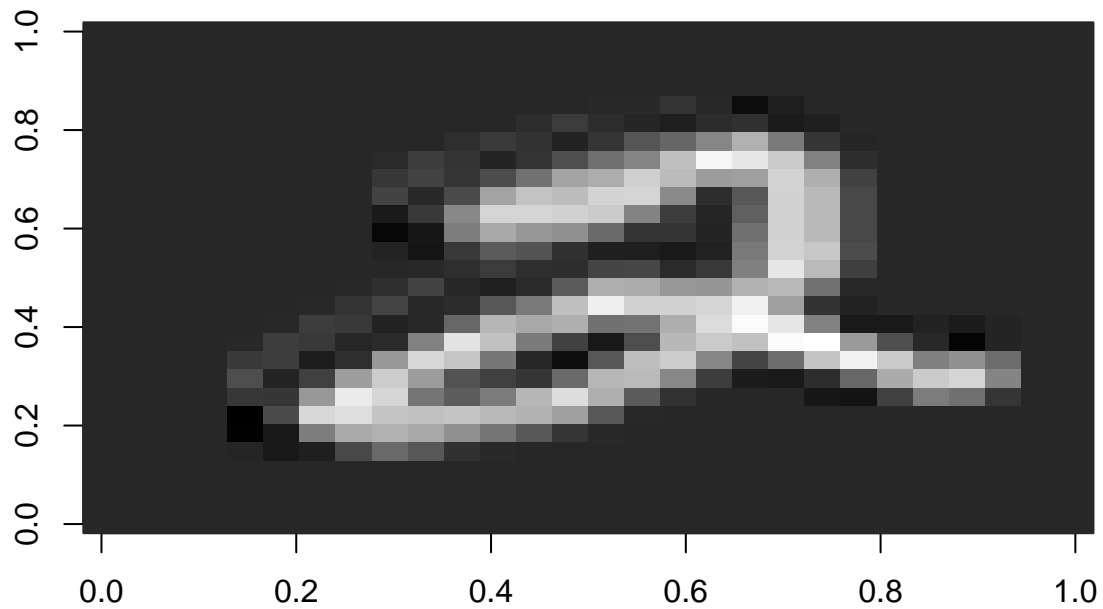
```
tf <- Sys.time()
print(tf - ti)
```

```
## Time difference of 0.009814024 secs
```

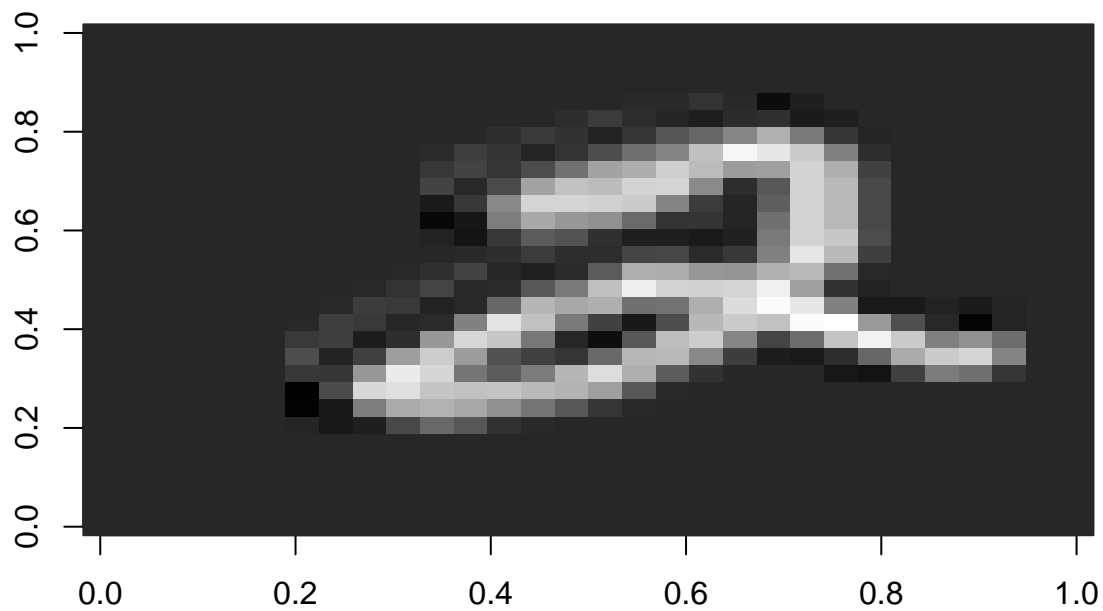
```
show <- 6
image(train_dat[["imgs"]][, , show], col = gray((0:255) / 255))
```



```
image(feature_maps_slow[, , show], col = gray((0:255) / 255))
```



```
image(feature_maps_fft[, , show], col = gray((0:255) / 255))
```



## Pooling

Below I program some pooling procedures to implement a pooling layer within the convolutional neural network.

## Backward Propagation