

Homework 2

1i) No learning is occurring in the neural network with weights initialized to zero because, relu will turn outputs to zero, the gradients of the weights will 0. In terms of back propagation, the gradients will all be 0, thus no learning / weight updates will occur.

1B) The large negative bias will close ReLU gates, when many neurons return 0, gradients can't flow during backpropagation, portions of the network can go dormant as closed relus can't change input parameters.

2ii) With Sigmoid the learning is a bit different compared to relu, instead of 0 learning occurring back propagation eventually the weights were learned. In terms of relu, instead of 0 learning occurring back propagation eventually the weights were learned. Similarly, it took almost the entire 4000 iterations and could only learn a simple linear boundary leaving a higher error compared to random weights which were able to learn a complex boundary with close to 0 test error. A boundary could still be learned because Sigmoid guarantees outputs between 0 and 1. Compared to random weights, learned because Sigmoid guarantees outputs between 0 and 1.

Home work 2

2 i)

$$(8 \cdot \underbrace{5 \cdot 5}_{\substack{\text{square} \\ \text{filter} \\ \text{size}}} \cdot 3) + 8$$

\uparrow \uparrow \uparrow
 # of filters filter size # of channels add bias terms per filter

608 parameters

2 ii) $W_{out} = \frac{W - F + \cancel{2P}}{S} + 1$ No zero padding

\nearrow

image is
square

$$\frac{W - F}{S} + 1 =$$

27 ~~$\frac{32 - 5}{1}$~~ + 1 = 28

28 x 28 x 3

3 i)

$$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0.5 \\ 0.5 & 0.25 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0.5 & 1 \\ 0.25 & 0.5 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0.5 & 0.25 \\ 1 & 0.5 \end{bmatrix}$$

3 ii)

$$\rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

~~They all have a~~

at least a single
each image has all 1 in
every quadrant if the 4×4
matrix is split up, every
quadrant will find a
max value of 1 resulting
in 4 identical pooling matrices.

problem C

Also in ~~pdf~~ Jupyter notebook

3C i)

1:	85.1233 %
2:	95.6817 %
3:	97.0417 %
4:	97.410 %
5:	98.0483 %
6:	98.3317 %
7:	98.5200 %
8:	98.6450 %
9:	98.8783 %
10:	98.9033 %

3C ii)

After A/B testing different layers and activation functions, I was able to get to a final test accuracy of 98.9% by using conv2d, linear and dropout layers and relu and max pooling activation functions, the model was fairly simple.

3C iii)

Yes, this is very computationally expensive, it can take a minute or two just to get a single epoch of training and 15-30 minutes to get the fully trained model, it's not an efficient way to hyperparameterize layers and activation functions.