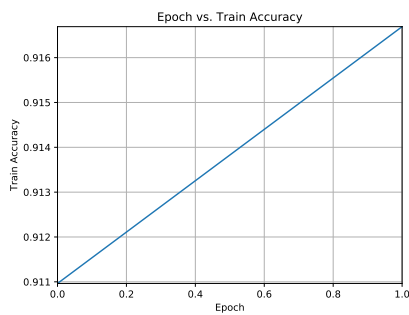


# 1 Report

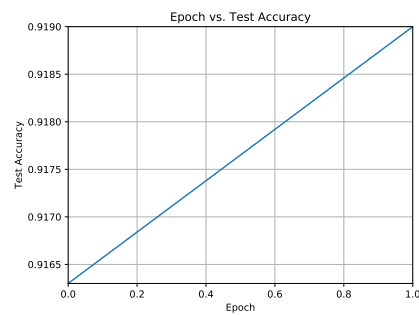
## 1.1 Part 1

I began by downloading TensorFlow's MNIST code for a simple neural network with one fully connected layer. After running the code as it was, I found an accuracy of 92% on the test set.

The results are shown in Figure ??.



(a) Part 1: Accuracy vs. Epoch as measured on the training set.



(b) Part 1: Accuracy vs. Epoch as measured on the test set.

## 1.2 Part 2

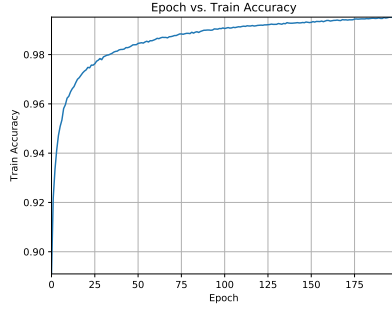
In the second part, I modified TensorFlow's MNIST code, as stated in the instructions, by expanding the neural network to include a total of 7 layers:

1. A 2d convolution layer with a  $10 \times 10$  kernel, a stride of 1, 32 output channels, and a ReLU activation function.
2. A max pooling layer with a  $2 \times 2$  kernel and a stride of 2.
3. Another 2d convolution layer with a  $5 \times 5$  kernel, a stride of 1, 16 output channels, and a ReLU activation function.
4. Another max pooling layer with a  $2 \times 2$  kernel and a stride of 2.
5. A flatten layer.
6. A fully connected layer with 1024 output channels and a ReLU activation function.

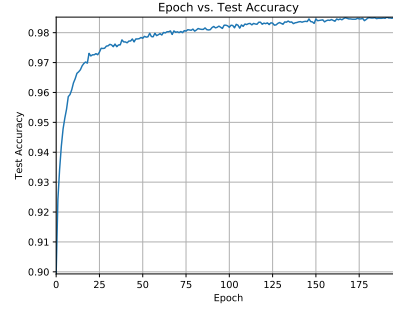
7. Another fully connected layer with 10 output channels and a Softmax activation function.

In addition, I used a cross-entropy loss function, the Adagrad optimizer with a learning rate of  $10^{-3}$ . Each of the weight variables was initialized by a truncated normal distribution with mean 0 and standard deviation 0.1, and each bias variables was initialized to 0.1.

The results are shown in Figure 3.



(a) Part 2: Accuracy vs. Epoch as measured on the training set.



(b) Part 2: Accuracy vs. Epoch as measured on the test set.

### 1.3 Part 3

In the third part, I implemented the neural network from Part 2 using only NumPy. To do this, I set up the class structure shown in Figure 3.

In order to compute each of the forward and backward steps, I implemented the equations in Section 2. In order to compute the convolutional and max pooling layers faster, I first reshaped the order 4 tensor into a matrix via  $x_{bijkl} \mapsto x_{(bIJ+iJ+j)(pQR+qR+r)}$  in order to use the full power of NumPy's vectorization capabilities. After taking a dot product, I would reshape the matrix back into a tensor. And in order to speed up the reshaping, I implemented it via a Cython module.

Despite my attempts to speed up the calculations, this implementation was far too slow to run 200 epochs, so I ran 5. The trial did not result in convergence, and I have not had sufficient time to work out all the bugs, although I am confident that I can given enough time.

## 2 Mathematical Methods

### 2.1 General

#### 2.1.1 Forward Propagation

$$\mathbf{x}^l = \sigma^{(l-1)}(\mathbf{z}^{(l-1)}) \quad (1)$$

$$\mathbf{z}^l = \mathbf{f}^l(\mathbf{x}^l, \mathbf{w}^l) + \mathbf{b}^l \quad (2)$$

#### 2.1.2 Backpropagation

$$\begin{aligned} \delta_{bi}^l &= -\eta \frac{\partial J_b}{\partial z_{bi}^l} = -\eta \frac{\partial J_b}{\partial x_{bj_0}^L} \frac{\partial x_{bj_0}^L}{\partial x_{bj_1}^{(L-1)}} \cdots \frac{\partial x_{j_{L-l-2}}^{(l+2)b}}{\partial x_{j_{L-l-1}}^{(l+1)b}} \frac{\partial x_{j_{L-l-1}}^{(l+1)b}}{\partial z_{bi}^l} \\ &= -\eta \frac{\partial J_b}{\partial x_{bj_0}^L} \left( \frac{\partial x_{bj_0}^L}{\partial z_{bk_0}^{(L-1)}} \frac{\partial z_{bk_0}^{(L-1)}}{\partial x_{bj_1}^{(L-1)}} \right) \cdots \left( \frac{\partial x_{j_{L-l-2}}^{(l+2)b}}{\partial z_{k_{L-l-2}}^{(l+1)b}} \frac{\partial z_{k_{L-l-2}}^{(l+1)b}}{\partial x_{j_{L-l-1}}^{(l+1)b}} \right) \frac{\partial x_{j_{L-l-1}}^{(l+1)b}}{\partial z_{bi}^l} \\ &= -\eta \frac{\partial J_b}{\partial x_{bj_0}^L} \frac{\partial \sigma_{bj_0}^{(L-1)}}{\partial z_{bk_0}^{(L-1)}} \frac{\partial f_{bk_0}^{(L-1)}}{\partial x_{bj_1}^{(L-1)}} \cdots \frac{\partial \sigma_{j_{L-l-2}}^{(l+1)b}}{\partial z_{k_{L-l-2}}^{(l+1)b}} \frac{\partial f_{k_{L-l-2}}^{(l+1)b}}{\partial x_{j_{L-l-1}}^{(l+1)b}} \frac{\partial \sigma_{bj_{L-l-1}}^l}{\partial z_{bi}^l} \end{aligned} \quad (3)$$

$$\delta_{bi}^{(L-1)} = -\eta \frac{\partial J_b}{\partial x_{bj}^L} \frac{\partial \sigma_{bj}^{(L-1)}}{\partial z_{bi}^{(L-1)}} \quad (4)$$

$$\delta_{bi}^{(l-1)} = \delta_{bj}^l \frac{\partial f_{bj}^l}{\partial x_{bk}^l} \frac{\partial \sigma_{bk}^{(l-1)}}{\partial z_{bi}^{(l-1)}} \quad (5)$$

$$\Delta b_{bi}^l = -\eta \frac{\partial J_b}{\partial b_{bi}^l} = -\eta \frac{\partial J_b}{\partial z_{bi}^l} = \delta_{bi}^l \quad (6)$$

$$\Delta w_{bij}^l = -\eta \frac{\partial J_b}{\partial w_{bij}^l} = -\eta \frac{\partial J_b}{\partial z_{bk}^l} \frac{\partial f_k^l}{\partial w_{bij}^l} = \delta_{bk}^l \frac{\partial f_k^l}{\partial w_{bij}^l} \quad (7)$$

$$\Delta b_i^l = \frac{1}{B} \sum_{b=0}^B \Delta b_{bi}^l \quad (8)$$

$$\Delta w_{ij}^l = \frac{1}{B} \sum_{b=0}^B \Delta w_{bij}^l \quad (9)$$

## 2.2 Layers

1. Flatten:

$$F(x_{bijk}^l) = x_{b(iJK+jK+k)}^l \quad (10)$$

$$\frac{\partial F_{bu}^l}{\partial x_{cv}^l} = \delta_{bc} \delta_{uv}, \quad \frac{\partial J}{\partial x_{c(iJK+jK+k)}^l} = (F_{cijk}^l)^{-1} \left( \frac{\partial J}{\partial z_{c(iJK+jK+k)}^l} \right) \quad (11)$$

2. Fully Connected (Dense):

$$D_{bj}^l(\mathbf{x}^l, \mathbf{w}^l) = \sum_i x_{bi}^l w_{ij}^l \quad (12)$$

$$\frac{\partial D_{bj}^l}{\partial x_{ck}^l} = \delta_{bc} w_{kj}^l, \quad \frac{\partial J}{\partial x_{ck}^l} = \frac{\partial J}{\partial z_{cj}^l} w_{kj}^l \quad (13)$$

$$\frac{\partial D_{bj}^l}{\partial w_{km}^l} = x_{bk}^l \delta_{jm}, \quad \frac{\partial J}{\partial w_{km}^l} = \frac{\partial J}{\partial z_{bm}^l} x_{bk}^l \quad (14)$$

3. 2D Convolution:

$$C_{bijk}^l(\mathbf{x}^l, \mathbf{w}^l) = \sum_{p,q,r} x_{b(s_1i+p)(s_2j+q)r}^l w_{(P-p)(Q-q)rk}^l \quad (15)$$

$$\frac{\partial C_{bijk}^l}{\partial x_{ctuv}^l} = \delta_{bc} w_{(P+s_1i-t)(Q+s_2j-u)vk}^l, \quad \frac{\partial J}{\partial x_{ctuv}^l} = \frac{\partial J}{\partial z_{cijk}^l} w_{(P+s_1i-t)(Q+s_2j-u)vk}^l \quad (16)$$

$$\frac{\partial C_{bijk}^l}{\partial w_{mnuv}^l} = x_{b(s_1i+P-m)(s_2j+Q-n)u}^l \delta_{kv}, \quad \frac{\partial J}{\partial z_{bijk}^l} x_{b(s_1i+P-m)(s_2j+Q-n)u}^l \quad (17)$$

4. Max Pooling  $2 \times 2$ :

$$P_{bijk}^l(\mathbf{x}) = \max_{p,q} x_{b(s_1i+p)(s_2j+q)k}^l \quad (18)$$

$$\frac{\partial P_{bijk}^l}{\partial x_{ctuv}^l} = \begin{cases} \delta_{bc} \delta_{it} \delta_{ju} \delta_{kv} & x_{bijk}^l = P(x_{bijk}^l) \\ 0 & \text{otherwise} \end{cases}, \quad \frac{\partial J}{\partial x_{ctuv}^l} = \begin{cases} \frac{\partial J}{\partial z_{bijk}^l} & x_{bijk}^l = P(x_{bijk}^l) \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

## 2.3 Activation Functions

1. ReLU:

$$R_{bi}^l(\mathbf{z}^l) = \begin{cases} z_{bi}^l & z_{bi}^l \geq 0 \\ 0 & z_{bi}^l < 0 \end{cases} \quad (20)$$

$$\frac{\partial R_{bi}^l}{\partial z_{cj}^l} = \begin{cases} \delta_{bc} \delta_{ij} & z_{ci}^l > 0 \\ \frac{1}{2} \delta_{bc} \delta_{ij} & z_{ci}^l = 0 \\ 0 & z_{ci}^l < 0 \end{cases}, \quad \frac{\partial J}{\partial z_{cj}^l} = \begin{cases} \frac{\partial J}{\partial x_{cj}^{l+1}} & z_{cj}^l > 0 \\ \frac{1}{2} \frac{\partial J}{\partial x_{cj}^{l+1}} & z_{cj}^l = 0 \\ 0 & z_{cj}^l < 0 \end{cases} \quad (21)$$

2. Softmax:

$$S_{bi}^l(\mathbf{z}^l) = \frac{e^{z_{bi}^l}}{\sum_k e^{z_{bk}^l}} \quad (22)$$

$$\frac{\partial S_{bi}^l}{\partial z_{cj}^l} = \begin{cases} \delta_{bc} S_{bj}^l (1 - S_{bj}^l) & i = j \\ -\delta_{bc} S_{bi}^l S_{bj}^l & i \neq j \end{cases}, \quad \frac{\partial J}{\partial z_{cj}^l} = \left( \frac{\partial J}{\partial x_{cj}^{l+1}} - \sum_i \frac{\partial J}{\partial x_{ci}^{l+1}} S_{ci}^l \right) S_{cj}^l \quad (23)$$

## 2.4 Cost Function

Softmax cross-entropy:

$$J(\mathbf{z}^{(L-1)}) = \frac{1}{B} \sum_{b=0}^B J_b(\mathbf{z}^{(L-1)}) \quad (24)$$

$$J_b(\mathbf{z}^{(L-1)}) = \sum_i -y_{bi} \log(z_{bi}^{(L-1)}) \quad (25)$$

$$\begin{aligned}
\frac{\partial J_b}{\partial z_{bj}^{(L-1)}} &= \frac{\partial J_b}{\partial x_{bk}^L} \frac{\partial S_k}{\partial z_{bj}^{(L-1)}} \\
&= -\frac{y_j}{x_{bj}^L} S_{bj}^{(L-1)} (1 - S_{bj}^{(L-1)}) + \sum_{k \neq j} \frac{y_k}{x_{bk}^L} S_{bk}^{(L-1)} S_{bj}^{(L-1)} \\
&= -\frac{y_j}{x_{bj}^L} x_{bj}^L (1 - x_{bj}^L) + \sum_{k \neq j} \frac{y_k}{x_{bk}^L} x_{bk}^L x_{bj}^L \\
&= -y_j (1 - x_{bj}^L) + \sum_{k \neq j} y_k x_{bj}^L \\
&= -y_j + x_{bj}^L \sum_k y_k \\
&= x_{bj}^L - y_j
\end{aligned} \tag{26}$$

## 2.5 This Case

### 2.5.1 Forward Propagation

$$L = 7 \tag{27}$$

$$\begin{aligned}
\sigma^0 &= \mathbf{R}(\mathbf{z}^0), & \mathbf{f}^0 &= \mathbf{C}(\mathbf{x}^0, \mathbf{w}^0) \\
\sigma^1 &= id(\mathbf{z}^1), & \mathbf{f}^1 &= \mathbf{P}(\mathbf{x}^1) \\
\sigma^2 &= \mathbf{R}(\mathbf{z}^2), & \mathbf{f}^2 &= \mathbf{C}(\mathbf{x}^2, \mathbf{w}^2) \\
\sigma^3 &= id(\mathbf{z}^3), & \mathbf{f}^3 &= \mathbf{P}(\mathbf{x}^3) \\
\sigma^4 &= id(\mathbf{z}^4), & \mathbf{f}^4 &= \mathbf{F}(\mathbf{x}^4) \\
\sigma^5 &= \mathbf{R}(\mathbf{z}^5), & \mathbf{f}^5 &= \mathbf{D}(\mathbf{x}^5, \mathbf{w}^5) \\
\sigma^6 &= \mathbf{S}(\mathbf{z}^6), & \mathbf{f}^6 &= \mathbf{D}(\mathbf{x}^6, \mathbf{w}^6)
\end{aligned} \tag{28}$$

## 2.5.2 Backpropagation

$$\begin{aligned}
\frac{\partial \sigma_{bi}^6}{\partial z_{bj}^6} &= \frac{\partial S_{bi}}{\partial z_{bj}^6}, & \frac{\partial f_{bi}^6}{\partial x_{bj}^6} &= \frac{\partial D_{bi}}{\partial x_{bj}^6}, & \frac{\partial f_{bi}^6}{\partial w_{bjk}^6} &= \frac{\partial D_{bi}}{\partial w_{bjk}^6} \\
\frac{\partial \sigma_{bi}^5}{\partial z_{bj}^5} &= \frac{\partial R_{bi}}{\partial z_{bj}^5}, & \frac{\partial f_{bi}^5}{\partial x_{bj}^5} &= \frac{\partial D_{bi}}{\partial x_{bj}^5}, & \frac{\partial f_{bi}^5}{\partial w_{bjk}^5} &= \frac{\partial D_{bi}}{\partial w_{bjk}^5} \\
\frac{\partial \sigma_{bi}^4}{\partial z_{bj}^4} &= \delta_{ij}, & \frac{\partial f_{bi}^4}{\partial x_{bj}^4} &= \frac{\partial F_{bi}}{\partial x_{bj}^4}, & & \\
\frac{\partial \sigma_{bi}^3}{\partial z_{bj}^3} &= \delta_{ij}, & \frac{\partial f_{bi}^3}{\partial x_{bj}^3} &= \frac{\partial P_{bi}}{\partial x_{bj}^3}, & & \\
\frac{\partial \sigma_{b p q r}^2}{\partial z^1} &= \frac{\partial R_{b i j k}}{\partial z_{b p q r}^2}, & \frac{\partial f_{b i j k}^2}{\partial x_{b p q r}^2} &= \frac{\partial C_{b i j k}}{\partial x_{b p q r}^2}, & \frac{\partial f_{b i j k}^2}{\partial w_{b p q r t}^2} &= \frac{\partial C_{b i j k}}{\partial w_{b p q r t}^2} \\
\frac{\partial \sigma_{bi}^1}{\partial z_{bj}^1} &= \delta_{ij}, & \frac{\partial f_{bi}^1}{\partial x_{bj}^1} &= \frac{\partial P_{bi}}{\partial x_{bj}^1}, & & \\
\frac{\partial \sigma_{b i j k}^0}{\partial z_{b p q r}^0} &= \frac{\partial R_{b i j k}}{\partial z_{b p q r}^0}, & \frac{\partial f_{b i j k}^0}{\partial x_{b p q r}^0} &= \frac{\partial C_{b i j k}}{\partial x_{b p q r}^0}, & \frac{\partial f_{b i j k}^0}{\partial w_{b p q r t}^0} &= \frac{\partial C_{b i j k}}{\partial w_{b p q r t}^0}
\end{aligned} \tag{29}$$

$$\begin{aligned}
\Delta b_{bi}^6 &= \delta_{bi}^6 = -\eta \frac{\partial J_b}{\partial z_{bi}^6} & \Delta w_{bij}^6 &= \delta_{bk}^6 \frac{\partial D_{bk}}{\partial w_{bij}^6} \\
\Delta b_{bi}^5 &= \delta_{bi}^5 = \delta_{bj}^6 \frac{\partial D_{bj}}{\partial x_{bk}^6} \frac{\partial R_{bk}}{\partial z_{bi}^5} & \Delta w_{bij}^5 &= \delta_{bk}^5 \frac{\partial D_{bk}}{\partial w_{bij}^5} \\
\delta_{bi}^4 &= \delta_{bj}^5 \frac{\partial D_{bj}}{\partial x_{bi}^5} \\
\delta_{b i j k}^3 &= \delta_{bp}^4 \frac{\partial F_{bp}}{\partial x_{b i j k}^4} \\
\Delta b_{b i j k}^2 &= \delta_{b i j k}^2 = \delta_{b p q r}^3 \frac{\partial P_{b p q r}}{\partial x_{b t u v}^3} \frac{\partial R_{b t u v}}{\partial z_{b i j k}^2} & \Delta w_{b i j k m}^2 &= \delta_{b p q r}^2 \frac{\partial C_{b p q r}}{\partial w_{b i j k m}^2} \\
\delta_{b i j k}^1 &= \delta_{b p q r}^2 \frac{\partial C_{b p q r}}{\partial x_{b i j k}^2} \\
\Delta b_{b i j k}^0 &= \delta_{b i j k}^0 = \delta_{b p q r}^1 \frac{\partial P_{b p q r}}{\partial x_{b t u v}^1} \frac{\partial R_{b t u v}}{\partial z_{b i j k}^0} & \Delta w_{b i j k m}^0 &= \delta_{b p q r}^0 \frac{\partial C_{b p q r}}{\partial w_{b i j k m}^0}
\end{aligned} \tag{30}$$

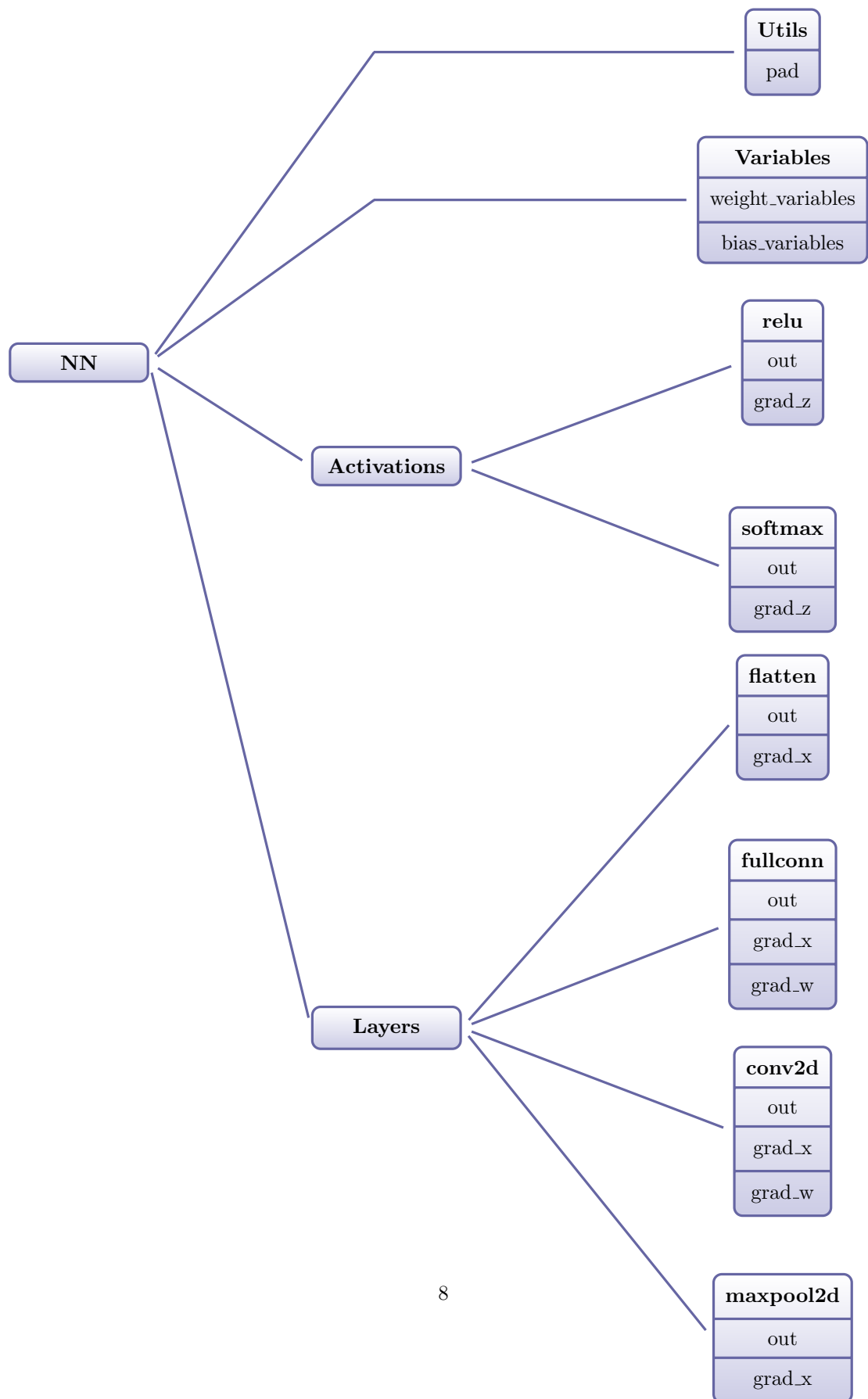


Figure 3: Class structure for neural network objects in Part 3.