

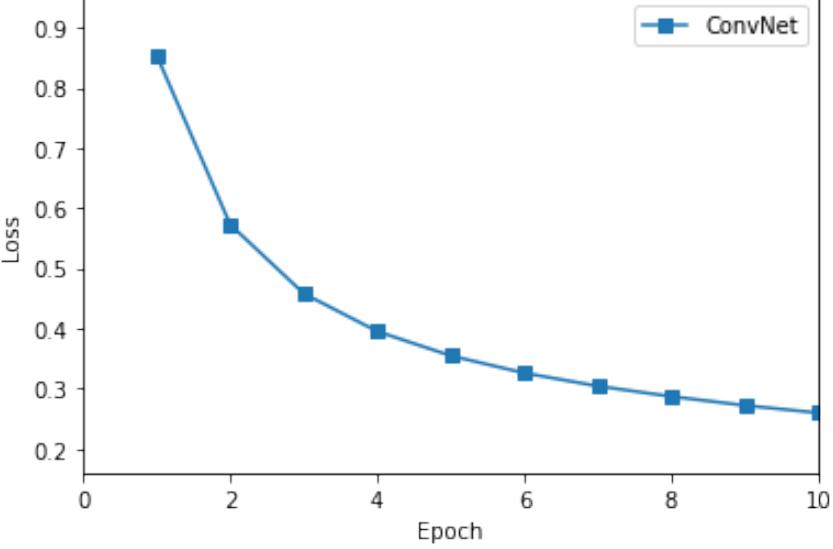
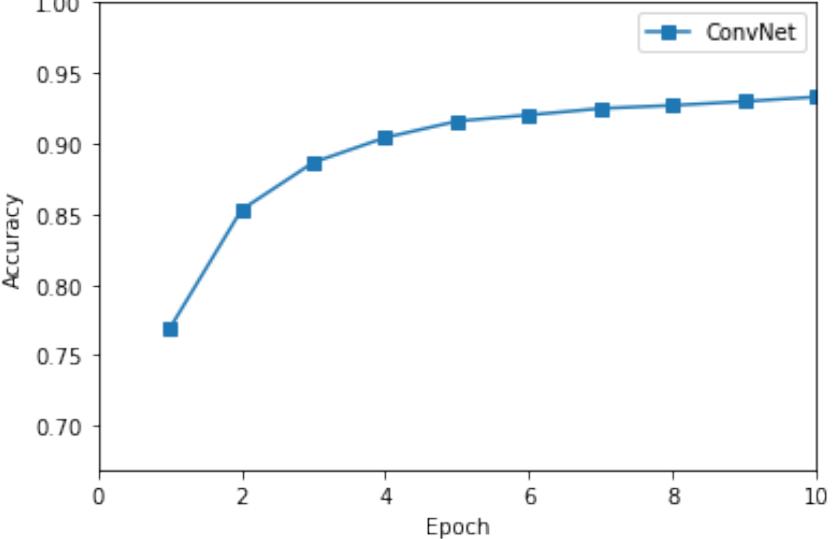
# HW2 – MNIST Classification with CNN

## 1. CNN for MNIST Classification

A CNN was trained on the MNIST dataset with the following hyper-parameters:

- batch\_size = 100
- max\_epoch = 10
- learning\_rate = 0.001
- weight\_decay = 0.005

The final training and test performance is recorded below

	Training Data	Test Data																						
Loss	0.3299  <table border="1"><thead><tr><th>Epoch</th><th>Loss</th></tr></thead><tbody><tr><td>1</td><td>0.85</td></tr><tr><td>2</td><td>0.58</td></tr><tr><td>3</td><td>0.46</td></tr><tr><td>4</td><td>0.40</td></tr><tr><td>5</td><td>0.36</td></tr><tr><td>6</td><td>0.33</td></tr><tr><td>7</td><td>0.31</td></tr><tr><td>8</td><td>0.29</td></tr><tr><td>9</td><td>0.28</td></tr><tr><td>10</td><td>0.26</td></tr></tbody></table>	Epoch	Loss	1	0.85	2	0.58	3	0.46	4	0.40	5	0.36	6	0.33	7	0.31	8	0.29	9	0.28	10	0.26	n/a
Epoch	Loss																							
1	0.85																							
2	0.58																							
3	0.46																							
4	0.40																							
5	0.36																							
6	0.33																							
7	0.31																							
8	0.29																							
9	0.28																							
10	0.26																							
Accuracy	0.9123  <table border="1"><thead><tr><th>Epoch</th><th>Accuracy</th></tr></thead><tbody><tr><td>1</td><td>0.77</td></tr><tr><td>2</td><td>0.85</td></tr><tr><td>3</td><td>0.89</td></tr><tr><td>4</td><td>0.91</td></tr><tr><td>5</td><td>0.92</td></tr><tr><td>6</td><td>0.925</td></tr><tr><td>7</td><td>0.93</td></tr><tr><td>8</td><td>0.935</td></tr><tr><td>9</td><td>0.94</td></tr><tr><td>10</td><td>0.94</td></tr></tbody></table>	Epoch	Accuracy	1	0.77	2	0.85	3	0.89	4	0.91	5	0.92	6	0.925	7	0.93	8	0.935	9	0.94	10	0.94	0.9197
Epoch	Accuracy																							
1	0.77																							
2	0.85																							
3	0.89																							
4	0.91																							
5	0.92																							
6	0.925																							
7	0.93																							
8	0.935																							
9	0.94																							
10	0.94																							

## 2. CNN vs MLP

Let us compare the performance between a 2-layer multilayer perceptron and the convolution neural network.

- The 2 layer MLP has the following architecture:
  - FC Layer → ReLU → FC Layer → ReLU → FC Layer → Softmax Cross-Entropy
- The CNN has the following architecture:
  - Conv → ReLU → MaxPool → Conv → ReLU → MaxPool → Reshape → FC Layer → ReLU → \*(Dropout) → FC Layer → Softmax Cross-Entropy

The hyper-parameters used for both models are listed below:

- batch\_size = 100
- max\_epoch = 10
- init\_std = 0.01
- learning\_rate\_SGD = 0.001
- weight\_decay = 0.005 (0 for dropout network)
- dropout\_rate = 0.2

### 2.1 MLP with Softmax Cross-Entropy and ReLU Activation

	Training	Validation	Test
Average Loss	0.7569	0.6797	n/a
Average Accuracy	0.8591	0.8988	0.8728

### 2.2 CNN with Softmax Cross-Entropy and ReLU Activation

	Training	Validation	Test
Average Loss	0.3299	0.2596	n/a
Average Accuracy	0.9123	0.9330	0.9197

### 2.3 CNN with Softmax Cross-Entropy and ReLU Activation and Dropout

	Training	Validation	Test
Average Loss	0.3399	0.2813	n/a
Average Accuracy	0.8956	0.9192	0.9305

In terms of training time, the CNN with dropout took the longest to complete 10 epochs, followed by the CNN without dropout and then the MLP. This makes sense because the CNN with dropout has the most complexity of the 3 architectures.

In terms of convergence, after 10 epochs the MLP is the only architecture that seems to have converged. Both the CNN with and without dropout have not reached convergence. This is probably a result of both the increased complexity of these models and the increased learning power.

In terms of accuracy, both CNNs have all around higher accuracy than the MLP model. Although the regular CNN had higher training and validation accuracies, it had a lower test accuracy compared to the CNN with dropout. This is due to the generalization cause by the dropout layer.

### 3. Dropout

	Training Data	Test Data																						
Loss	<p>0 . 3399</p> <table border="1"> <caption>Data for Loss vs Epoch (ConvNet)</caption> <thead> <tr> <th>Epoch</th> <th>Loss</th> </tr> </thead> <tbody> <tr><td>1</td><td>0.70</td></tr> <tr><td>2</td><td>0.52</td></tr> <tr><td>3</td><td>0.43</td></tr> <tr><td>4</td><td>0.39</td></tr> <tr><td>5</td><td>0.36</td></tr> <tr><td>6</td><td>0.34</td></tr> <tr><td>7</td><td>0.31</td></tr> <tr><td>8</td><td>0.30</td></tr> <tr><td>9</td><td>0.29</td></tr> <tr><td>10</td><td>0.29</td></tr> </tbody> </table>	Epoch	Loss	1	0.70	2	0.52	3	0.43	4	0.39	5	0.36	6	0.34	7	0.31	8	0.30	9	0.29	10	0.29	n/a
Epoch	Loss																							
1	0.70																							
2	0.52																							
3	0.43																							
4	0.39																							
5	0.36																							
6	0.34																							
7	0.31																							
8	0.30																							
9	0.29																							
10	0.29																							
Accuracy	<p>0 . 8956</p> <table border="1"> <caption>Data for Accuracy vs Epoch (ConvNet)</caption> <thead> <tr> <th>Epoch</th> <th>Accuracy</th> </tr> </thead> <tbody> <tr><td>1</td><td>0.78</td></tr> <tr><td>2</td><td>0.84</td></tr> <tr><td>3</td><td>0.87</td></tr> <tr><td>4</td><td>0.89</td></tr> <tr><td>5</td><td>0.90</td></tr> <tr><td>6</td><td>0.90</td></tr> <tr><td>7</td><td>0.91</td></tr> <tr><td>8</td><td>0.91</td></tr> <tr><td>9</td><td>0.92</td></tr> <tr><td>10</td><td>0.92</td></tr> </tbody> </table>	Epoch	Accuracy	1	0.78	2	0.84	3	0.87	4	0.89	5	0.90	6	0.90	7	0.91	8	0.91	9	0.92	10	0.92	0 . 9305
Epoch	Accuracy																							
1	0.78																							
2	0.84																							
3	0.87																							
4	0.89																							
5	0.90																							
6	0.90																							
7	0.91																							
8	0.91																							
9	0.92																							
10	0.92																							

The dropout layer acts as a regularizer. I have tested multiple locations as well as different dropout rates and have observed the following:

- Dropout layers placed after Conv → ReLU segments negatively affect classification performance. However, if the dropout rate is very small (<0.1) at these layers it results in comparable performance with no dropout.
- Dropout layers placed after Linear → ReLU segments are the most effective at regularizing the model
- Dropout rates that are higher than 0.3 cause the model to underperform. This may be because CNN with dropout takes longer to converge, and thus must be trained over more epochs

## 4. Appendix: Mathematical Derivations

### Batch CNN

BATCH CONVOLUTION LAYER

$$\left\{ \begin{array}{l} X \in \mathbb{R}^{B \times F_l \times H \times W} \\ Z \in \mathbb{R}^{B \times F_{l+1} \times (H-d_1+1) \times (W-d_2+1)} \\ W \in \mathbb{R}^{F_l \times F_{l+1} \times d_1 \times d_2} \end{array} \right.$$

B: batch size  
 $F_l$ : #feature maps at layer  $l$   
 H: height of feature map  
 W: width " "  
 $d_1$ : height of filter  
 $d_2$ : width of filter

FORWARD PASS

$$\begin{aligned}
 Z_{m,n}^{(N,Q)} &= \sum_P \left( \sum_{i,j} X_{m+i,n+j}^{(N,P)} W_{i,j}^{(P,Q)} \right) + b^{(Q)} \\
 &= \underbrace{\sum_{P,i,j} X_{m+i,n+j}^{(N,P)} W_{i,j}^{(P,Q)}}_{\text{equal-valid correlation on } P} + b^{(Q)}
 \end{aligned}$$

valid correlation on  $i, j$

$$\Rightarrow Z^{(N,Q)} = X \circ_{\text{valid}} [W^T]^{(Q)} + b^{(Q)}$$

\* where  $T(x,y)$  is the tensor transpose between axis x and axis y.

BACKWARD PASS

$$\begin{aligned}
 \frac{\partial L}{\partial X_{l,k}^{(N,P)}} &= \sum_{Q,M,N} \frac{\partial L}{\partial Z_{m,n}^{(N,Q)}} \cdot \frac{\partial Z_{m,n}^{(N,Q)}}{\partial X_{l,k}^{(N,P)}} \\
 &= \sum_{Q,M,N} \delta_{M,N}^{(N,Q)} W_{l-M, k-N}^{(P,Q)}
 \end{aligned}$$

equal-valid correlation on Q  
 full convolution on M,N

$$\Rightarrow \frac{\partial L}{\partial X^{(N,P)}} = \delta^{(N)} \circ_{\text{valid}(o)} W^{(P)}$$

\* full(1,2)

$$= \delta^{(N)} *_{\substack{\text{valid}(o) \\ \text{full}(1,2)}} \text{flip}_x [W^{(P)}]$$

\* where  $\text{flip}_x$  denotes flipping the tensor along axis x.

The gradients of the weights and biases are...

$$\frac{\partial L}{\partial W_{i,j}^{(P,Q)}} = \sum_{N,M,n} \frac{\partial L}{\partial Z_{m,n}^{(N,Q)}} \cdot \frac{\partial Z_{m,n}^{(N,Q)}}{\partial W_{i,j}^{(P,Q)}}$$

$$= \underbrace{\sum_{N,M,n} \delta_{m,n}^{(N,Q)} \cdot X_{i+m,j+n}^{(N,P)}}_{\substack{\text{equal-valid correlation on } N \\ \text{valid correlation on } m,n}}$$

$$\Rightarrow \frac{\partial L}{\partial W^{(P,Q)}} = [X^T \circ_{\text{valid}} \delta^{(Q)}]^{(P)}$$

$$\frac{\partial L}{\partial b^{(Q)}} = \sum_{N,M,n} \frac{\partial L}{\partial Z_{m,n}^{(N,Q)}} \cdot \frac{\partial Z_{m,n}^{(N,Q)}}{\partial b^{(Q)}}$$

$$= \sum_{N,M,n} \delta_{m,n}^{(N,Q)} \cdot 1$$

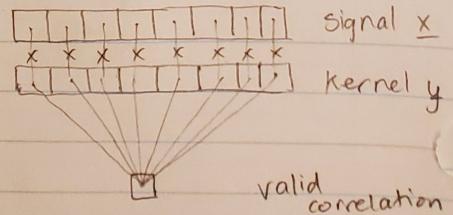
$$\Rightarrow \frac{\partial L}{\partial b^{(Q)}} = \sum_{N,M,n} \delta_{m,n}^{(N,Q)} \quad //$$

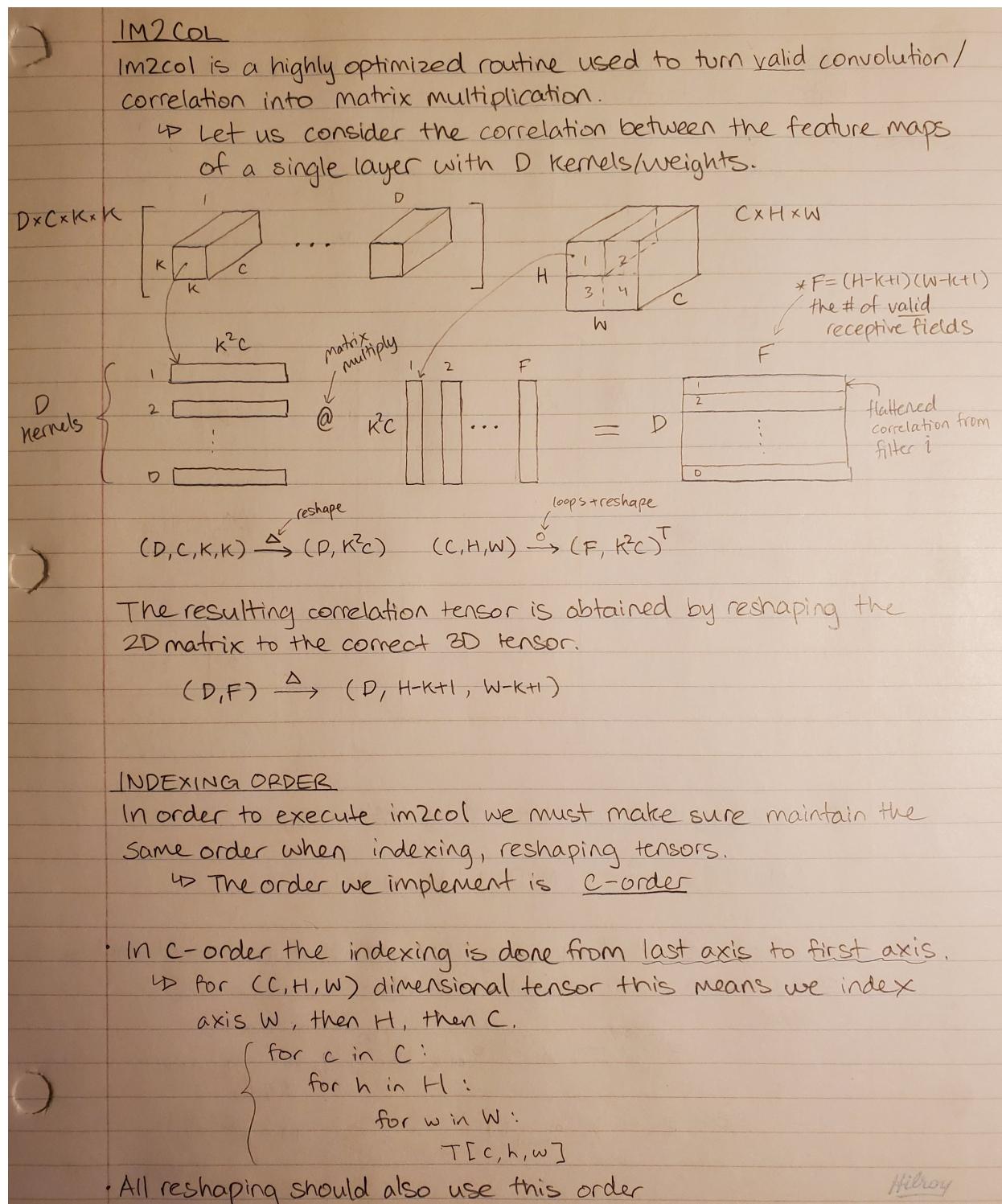
### EQUAL-VALID CORRELATION

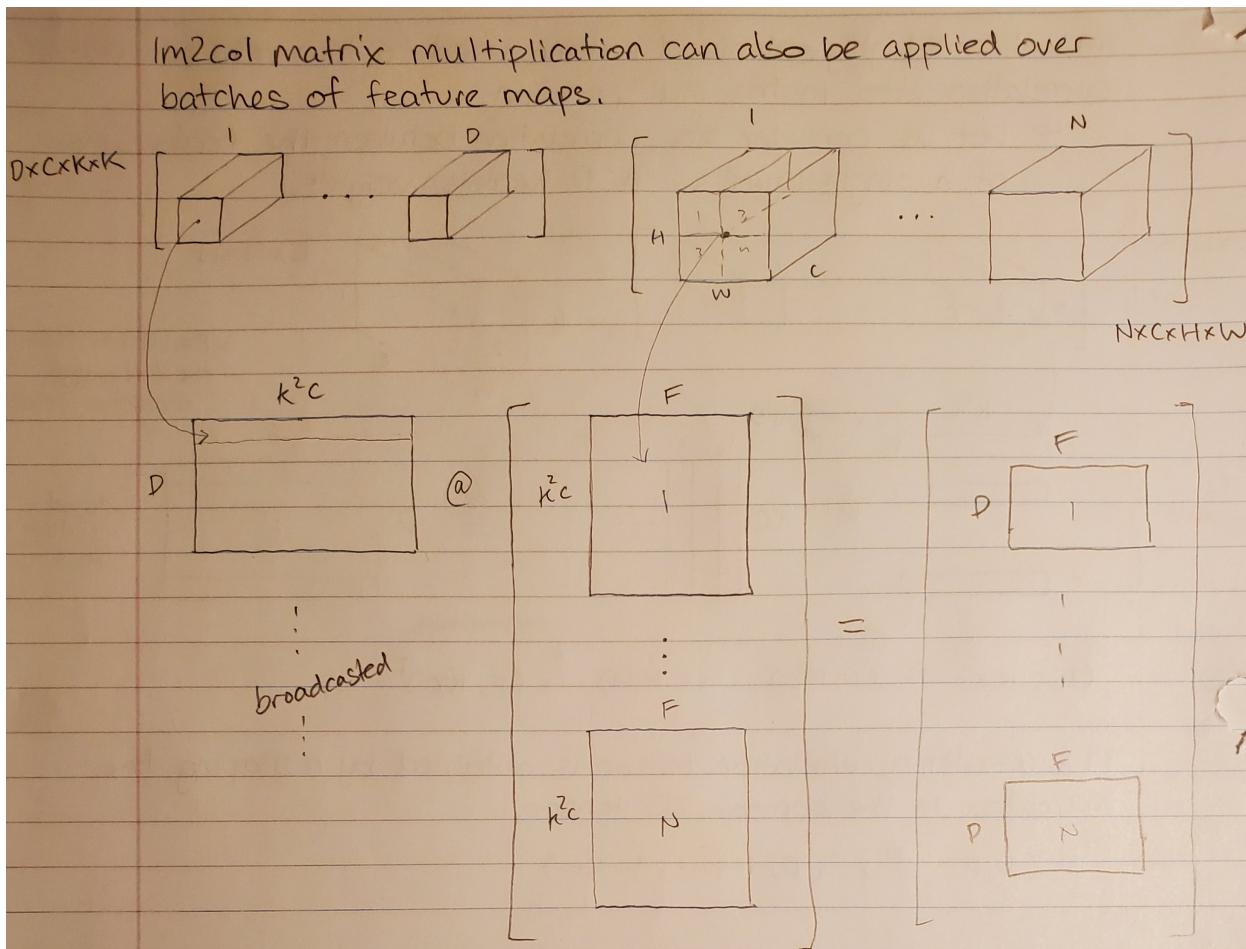
We can show that the dot product of 2 vectors is actually a form of valid correlation.

$\hookrightarrow x \in \mathbb{R}^d$     } equal dimensions!  
 $\hookrightarrow y \in \mathbb{R}^d$

$$x^T y = \sum_i x_i y_i = x \circ_{\text{valid}} y$$



**im2col**



The resulting tensor is obtained by reshaping...

$$(N, D, F) \xrightarrow{\Delta} (N, D, H-K+1, W-K+1)$$

- Therefore instead of applying 3D correlation in forward/backward pass we can apply im2col correlation

Ex.

$$Z^{(N,D)} = X^{(N)} \circ_{\text{valid}} W^{(D)}$$

$$\Rightarrow Z = X \circ_{\text{im2col}} W$$

$X \in \mathbb{R}^{N, C, H, W}$   
 $W \in \mathbb{R}^{D, C, K, K}$   
 $Z \in \mathbb{R}^{N, D, H-K+1, W-K+1}$