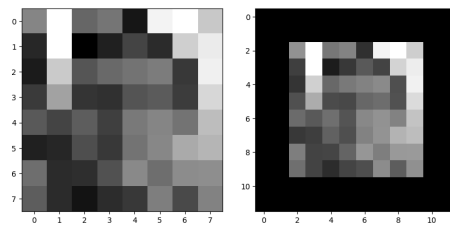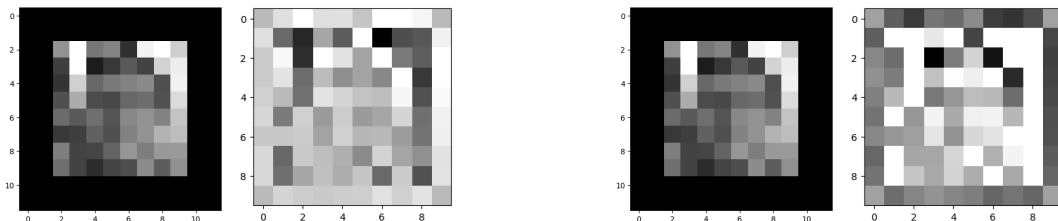**Padding**: We want the kernel (filter) to move over the image matrix (the next step). However, since we want it to have an equal effect on every value, **we fill the edges of our matrix with '0'.**
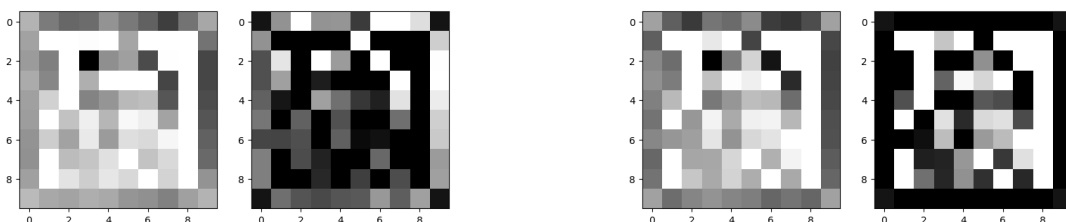


**Full Correlation**: The image matrix is traversed with a special kernel (filter). The values of the kernel and the corresponding image area are multiplied, then the matrix values are summed up and turned into a single float value, which is then written to the new area. **This is used to learn features in the image (edges, corners, patterns).**
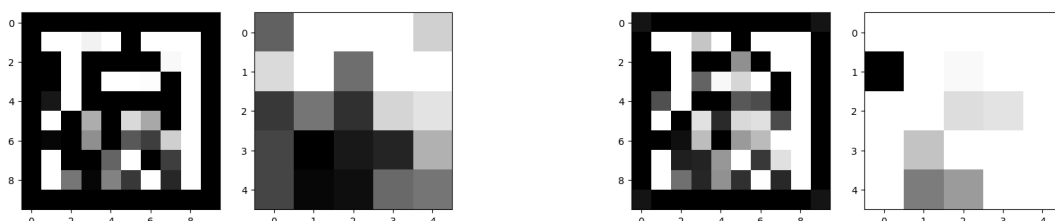


**Activation (ReLU, Sigmoid):**

- **ReLU: It works by setting negative values to zero and keeping positive values unchanged.** This function helps gradients propagate more effectively, especially in deep networks.
- **Sigmoid**: It provides an S-curve activation with **outputs in the range [0, 1].** It is particularly used in binary classification problems.



**Pooling (Max Pooling):** We move a 3x3 matrix over the image. The maximum value in the area covered by this 3x3 matrix is taken and written to the new matrix. **This is used to reduce the size of the image.**

**Initialization of Filters and Weights:**

*We have two filters, and each filter has a size of 5x5:*

$f1$ = **initializeFilter**(size, scale = 1.0)
$f2$ = **initializeFilter**(size, scale = 1.0)

*And We've got 2 hidden layers, each layer has 10 weight. They are size of 10x500 and 10x10:*

$w3$ = **initializeWeight**(size)
$w4$ = **initializeWeight**(size)

| implementation of **initializeFilter:** | implementation of **initializeWeight:** |
|---|---|
| $stddev = \dfrac{scale}{\sqrt{\prod\limits_{k=0}^{n} size[k]}}$  $\mathbb{N}(0, stddev)$ | $\mathbb{N}(0, 0.01)$ |

> $\mathbb{N}(\mu, \sigma)$: *denotes a normal distribution with mean $\mu$ and variance $\sigma$.*

**Randomizing Images to Prevent Over-fitting:**

1. **Scaling:**

   Randomly select a **scale factor between 0.9 and 1.1**, **resize the input image** with selected scale factor.

2. **Rotation:**

   Randomly **selects an angle between -45 and 45 degrees**, calculate center of the image then **apply rotation to the scaled image** using the rotation matrix.

3. **Translation:**

   Randomly select **translations in the x and y directions between -2 and 2 pixels**, construct a translation matrix based on the selected translations. **apply translation to the rotated image** using the translation matrix.

4. **Additive Noise:**

   Randomly select a **probability between 0 and 0.3**, generate a binary mask of the same shape as the translated image based on the selected probability.

   **Create random noise values of the same shape as the translated image**, scaled between -0.5 and 0.5, **combine the translated image with the noise values based on the binary mask.**

   **Clip the resulting noisy image to ensure pixel values are within the range [0, 1].**

**Forward-Propagation:**

1. **First Convolution Layer (Conv2D + ReLU):**

    $conv1$ = **Conv2D**(*image*, f1, b1, conv_stride)
    $conv1$ = **ReLU**(conv1) = **max**(0, conv1)

2. **Second Convolution Layer (Conv2D + ReLU):**

    $conv2$ = **Conv2D**(*conv1*, f2, b2, conv_stride)
    $conv2$ = **ReLU**(conv2) = **max**(0, conv2)

3. **Max Pooling:**

    $pooled$ = **MaxPool2D**(*conv2*, pool_size, pool_stride)

4. **Flattening the Pooled Layer:**

    $fc$ = **Flatten**(*pooled*) = pooled**.reshape**((shape of *pooled* in 1 dimension))

5. **Fully Connected Layer:**

    $z$ = W3 · $fc$ + b3
    $z$ = **ReLU**($z$) = **max**(0, $z$)
    $out$ = W4 · $z$ + b4
    $probs$ = **Softmax**(*out*)

**Calculating Cross-Entropy Loss:**

> *probs*: The result from the final layer of the Fully Connected Layer (Softmax), which is our prediction. **($\hat{Y}$)**

> *label:* The true label **(Y)**

> $\epsilon$: A very small value added to avoid taking the log(10). **(for exp. $10^{-10}$)**

$probs$ = **max**($\epsilon$, **min**(*probs*, 1−$\epsilon$))
$loss$ = −$\sum$(label · log(*probs*))

**Backpropagation:**

1. **Compute** *dout*:

   > *probs*: The result from the final layer of the Fully Connected Layer (Softmax), which represents our prediction. **(Ŷ)**

   > *label:* The true label. **(Y)**

   $dout = probs - \text{label}$

2. **Compute** *dw4* (Partial Derivative of *Loss* with respect to $W_4$):

   $dw4 = \frac{\partial Loss}{\partial W4} = \frac{\partial Loss}{\partial out} \cdot \frac{\partial out}{\partial W4} = dout \cdot z^{\mathrm{T}}$

3. **Compute** *db4* (Partial Derivative of *Loss* with respect to $b_4$):

   $db4 = \frac{\partial Loss}{\partial b4} = \frac{\partial Loss}{\partial out} \cdot \frac{\partial out}{\partial b4} = \sum dout$

4. **Compute** *dz* and *ReLU* (Partial Derivative of *Loss* with respect to *z*):

   $dz = \frac{\partial Loss}{\partial z} = \frac{\partial Loss}{\partial out} \cdot \frac{\partial out}{\partial z} = W_4^T \cdot dout$

   $dz = \text{ReLU}(z) \cdot dz$

5. **Compute** *dw3* (Partial Derivative of *Loss* with respect to $W_3$):

   $dW3 = \frac{\partial Loss}{\partial W3} = \frac{\partial Loss}{\partial z} \cdot \frac{\partial z}{\partial W1} = dz \cdot fc^{\mathrm{T}}$

6. **Compute** *db3* (Partial Derivative of *Loss* with respect to $b_3$):

   $db3 = \frac{\partial Loss}{\partial b3} = \frac{\partial Loss}{\partial z} \cdot \frac{\partial z}{\partial b3} = \sum dz$

7. **Compute** *dFc* and *dPool,* Then flatten *dPool*:

   $dFc = \frac{\partial Loss}{\partial Fc} = W_3^T \cdot dz$

   $dPool = dFc.\textbf{reshape}((\text{shape of } pooled))$

8. **Compute** *dConv2* and Backpropagate Max Pooling Layer:

   $dConv2 = \textbf{MaxPool2D\_Backward}(dPool, conv2, \text{pool\_size, pool\_stride})$

   $dConv2 = \textbf{ReLU}(conv2) \cdot dConv2$

9. **Compute** *dConv1, dF2, db2* and Backpropagate Convolution Layer:

   $dConv1, df2, db2 = \textbf{Convolution\_Backward}(dConv2, conv1, f2, \text{conv\_stride})$

   $dConv1 = \textbf{ReLU}(conv1) \cdot dConv1$

10. **Compute** *dImage, dF1, db1* and Backpropagate Convolution Layer:

    $dImage, df1, db1 = \textbf{Convolution\_Backward}(dConv, conv2, f1, \text{conv\_stride})$

**Adam (Adaptive Moment Estimation) Optimization:**
It adapts the learning rate for each parameter by computing individual learning rates for each parameter from estimates of first and second moments of the gradients.

1. **Initialize parameters**
   $\theta$, (Model Parameters)
   $m = 0$, (First Moment)
   $v = 0$, (Second Moment)

2. For each ***parameter ($\theta$)*** in ***Parameters***:
   > ***g*** *is our gradient such as df1, db1, dw3*
   > $\epsilon$ *is a very small number such as $10^{-10}$, that prevents us from getting dived by*
   *0*

   $$m = \beta_1 \cdot m + (1-\beta_1) \cdot g / batch\_size$$
   $$v = \beta_2 \cdot v + (1-\beta_2) \cdot (g / batch\_size)^2$$
   $$\theta = \theta - \frac{learningRate \cdot m}{\sqrt{v} + \epsilon}$$

**Using Multiple Threads For Faster Computing:**
Due to the extensive size of our training dataset, it's beneficial to split our batches into smaller sections. This enables us to leverage asynchronous processing with multiple threads. By doing so, we can distribute the workload across various threads, allowing for simultaneous computation on smaller subsets of data.
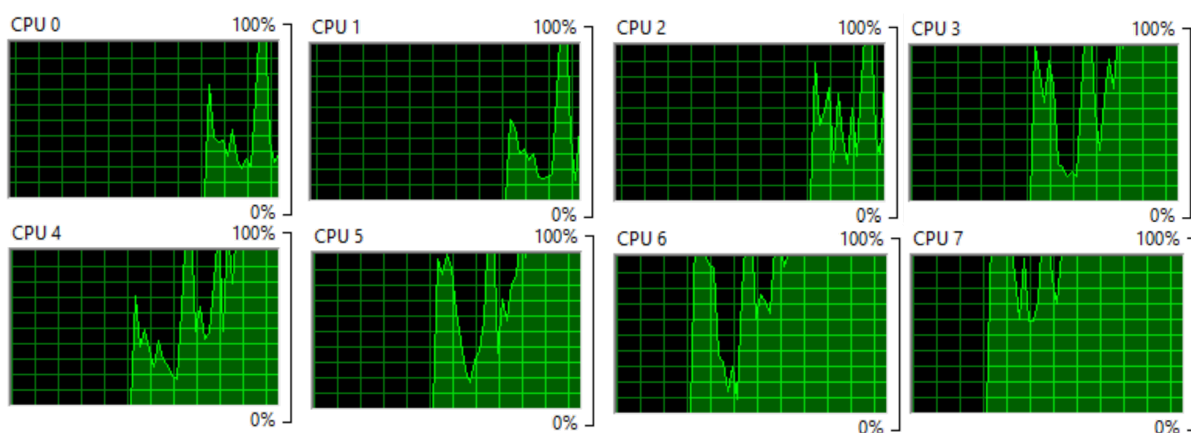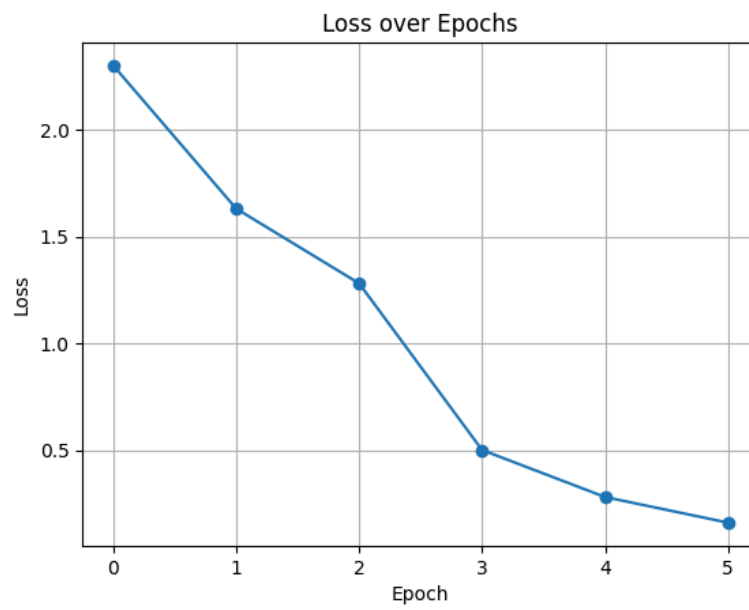
Batch Size = 1024
Total Training Data = 81.000
$$\text{Number of Epochs} = \frac{Total\ number\ of\ training\ data}{BatchSize}$$
Number of Processes = 8

So we split our batches into 8 sections, then asynchronously processing data with multiple threads.

**Loss Values Over Epochs:**



*with final accuracy of **93.252%***

**T.L = True Label, P = Prediction**



| | | |
|---|---|---|
| **T.L : 0,  P : 0** | **T.L : 5,  P : 5** | **T.L : 2,  P : 7** |
| **T.L : 9,  P : 9** | **T.L : 6,  P : 0** | **T.L : 7,  P : 7** |