# Demystifying the Mathematics Behind Convolutional Neural Networks (CNNs)

## Overview

- Convolutional neural networks (CNNs) are all the rage in the deep learning and computer vision community
- How does this CNN architecture work? We'll explore the math behind the building blocks of a convolutional neural network
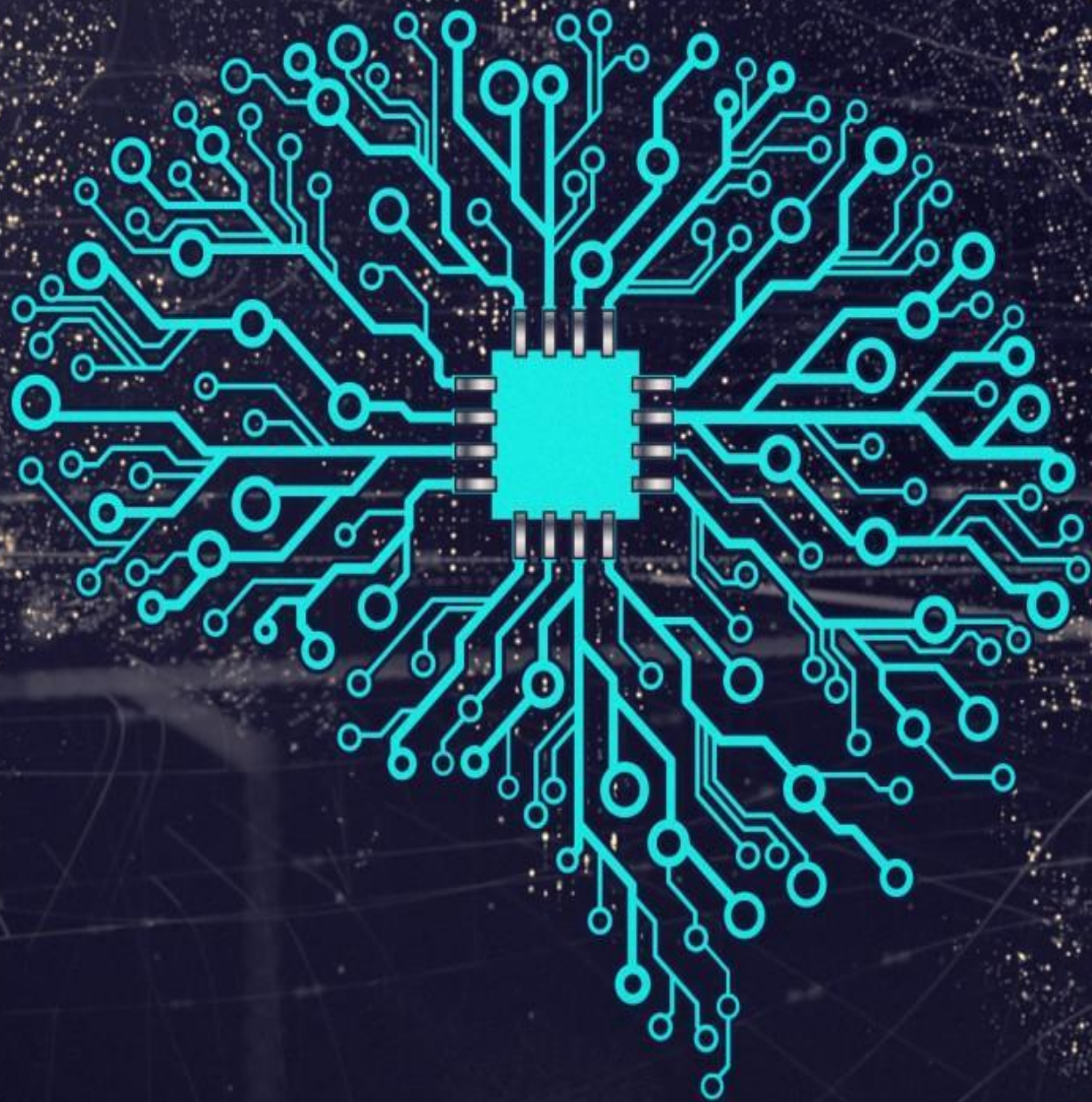- We will also build your own CNN from scratch using NumPy

## Introduction

Convolutional neural network (CNN) – almost sounds like an amalgamation of biology, art and mathematics. In a way, that's exactly what it is (and what this article will cover).

CNN-powered deep learning models are now ubiquitous and you'll find them sprinkled into various computer vision applications across the globe. Just like XGBoost and other popular machine learning algorithms, convolutional neural networks came into the public consciousness through a hackathon (the ImageNet competition in 2012).

These neural networks have caught inspiration like fire since then, expanding into various research areas. Here are just a few popular computer vision applications where CNNs are used:

- Facial recognition systems
- Analyzing and parsing through documents
- Smart cities (traffic cameras, for example)
- Recommendation systems, among other use cases

But why does a convolutional neural network work so well? How does it perform better than the traditional ANNs (Artificial neural network)? Why do deep learning experts love it?

To answer these questions, we must understand how a CNN actually works under the hood. In this article, we will go through the mathematics behind a CNN model and we'll then build our own CNN from scratch.

If you prefer a course format where we cover this content in stages, you can enrol in this free course:
[Convolutional Neural Networks from Scratch](#)

*Note: If you're new to neural networks, I highly recommend checking out our popular free course:*

- [*Introduction to Neural Networks*](#)

# Table of Contents

2. Back ward Propagation

# Introduction to Neural Networks

Neu ral  Network s  are at th e  core   of  all  <u>deep  learning  algorithms</u> . Bu t  before  you  deep  dive  in to  th ese  algorith ms, it's impor tan t to h ave a good u n derstan din g of th e con cept of <u>neural networks</u>.

T h ese n eu ral n etwork s try to  mimic  th e h u man brain  an d its  learn in g process. Lik e a brain  tak es th e in pu t, processes it an d gen erates some ou tpu t, so does th e n eu ral n etwork .

T h ese th ree   action s – **r ecei v i ng i nput, pr ocessi ng i nf or mati on, gener ati ng output**  – are represen ted  in th e form  of  layers in a n eu ral n etwork – in pu t, h idden  an d ou tpu t. Below  is a sk eleton  of wh at a n eu ral n etwork look s lik e:



T h ese in dividu al  u n its  in th e layers  are called  **neur ons**. T h e complete  train in g process of a n eu ral n etwork in volves two steps.

# 1. Forward Propagation

Images are fed in to the input layer in the form of numbers. These numerical values denote the intensity of pixels in the image. Then neurons in the hidden layers apply a few mathematical operations on these values (which we will discuss later in this article).

In order to perform these mathematical operations, there are certain parameter values that are randomly initialized. Post these mathematical operations at the hidden layer, the result is sent to the output layer which generates the final prediction.

# 2. Backward Propagation

Once the output is generated, the next step is to compare the output with the actual value. Based on the final output, and how close or far this is from the actual value (error), the values of the parameters are updated. The forward propagation process is repeated using the updated parameter values and new outputs are generated.

This is the base of any neural network algorithm. In this article, we will look at the forward and backward propagation steps for a convolutional neural network!

# Convolutional Neural Network (CNN) Architecture

Consider this – you are asked to identify objects in two given images. How would you go about doing that? Typically, you would observe the image, try to identify different features, shapes and edges from the image. Based on the information you gather, you would say that the object is a dog or a car and so on.

This is precisely what the hidden layers in a <u>C N N</u> do – find features in the image. The convolutional neural network can be broken down into two parts:

- **The convolution layers**: Extracts features from the input
- **The fully connected (dense) layers**: Uses data from convolution layer to generate output



As we discussed in the previous section, there are two important processes involved in the training of any neural network:

1. **Forward Propagation:** Receive input data, process the information, and generate output
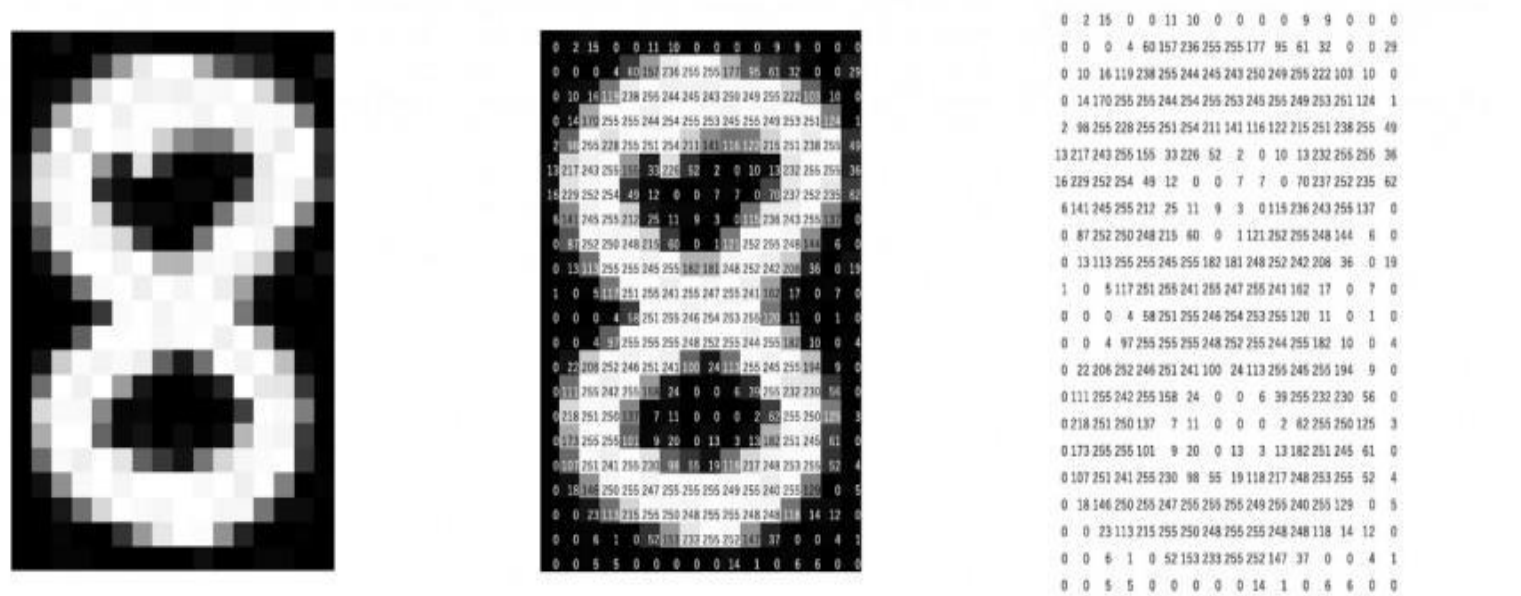2. **Backward Propagation:** Calculate error and update the parameters of the network

We will cover both of these one by one. Let us start with the forward propagation process.

# Convolutional Neural Network (CNN): Forward Propagation

## Convolution Layer

You know how we look at images and identify the object's shape and edges? A convolutional neural network does this by comparing the pixel values.

Below is an image of the number 8 and the pixel values for this image. Take a look at the image closely. You would notice that there is a significant difference between the pixel values around the edges of the number. Hence, a simple way to identify the edges is to compare the neighboring pixel value.



Do we need to traverse pixel by pixel and compare these values? No! To capture this information, the image is convolved with a filter (also known as a 'kernel').

Convolution is often represented mathematically with an asterisk * sign. If we have an input image represented as X and a filter represented with $f$, then the expression would be:

$$Z = X * f$$

*Note: To learn how filters capture information about the edges, you can go through this article:*

- *[Beginner-Friendly Techniques to Extract Features from Image Data](#)*

Let us understand the process of convolution using a simple example. Consider that we have an image of size 3 x 3 and a filter of size 2 x 2:

|    |    |    |
|----|----|----|
| 1  | 7  | 2  |
| 11 | 1  | 23 |
| 2  | 2  | 2  |

|   |   |
|---|---|
| 1 | 1 |
| 0 | 1 |

T h e filter  goes th rou gh  th e patch es of images, performs an elemen t-wise mu ltiplication , an d th e valu es are su mmed u p:

```
(1x1 + 7x1 + 11x0 + 1x1) = 9 (7x1 + 2x1 + 1x0 + 23x1)  = 32 (11x1  + 1x1 + 2x0 + 2x1) = 14 (1x1 + 23x1 + 2x0 +
2x1) = 26
```



Look  at th at closely – you 'll  n otice  th at th e filter  is  con siderin g a small  por tion  of th e image  at a time.  W e can also imagin e th is as a sin gle image brok en  down  in to  smaller patch es, each of wh ich  is con volved  with th e filter.



I n th e above example,  we   h ad an in pu t of sh ape (3,  3) an d a filter  of sh ape (2,  2). Sin ce  th e dimen sion s of image  an d filter  are  very  small,  it's  easy  to  in terpret th at th at th e  sh ape of  th e  ou tpu t matrix  is  (2,  2).  Bu t h ow wou ld  we  fin d  th e  sh ape of  an ou tpu t  for  more complex  in pu ts or  filter  dimen sion s? T h ere  is  a simple formu la to do  so:
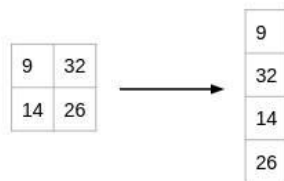
```
Dimension  of image = (n, n) Dimension  of filter  = (f,f) Dimension  of output  will be ((n-f+1)  , (n-f+1))
```

You sh ou ld h ave a good u n derstan din g of h ow a con volu tion al layer work s at th is poin t. Let u s move  to th e n ext  par t of th e  C N N  arch itectu re.

# Fully Connected Layer

So far, the convolution layer has extracted some valuable features from the data. These features are sent to the fully connected layer that generates the final results. The fully connected layer in a CNN is nothing but the traditional neural network!

The output from the convolution layer was a 2D matrix. Ideally, we would want each row to represent a single input image. In fact, the fully connected layer can only work with 1D data. Hence, the values generated from the previous operation are first converted into a 1D format.
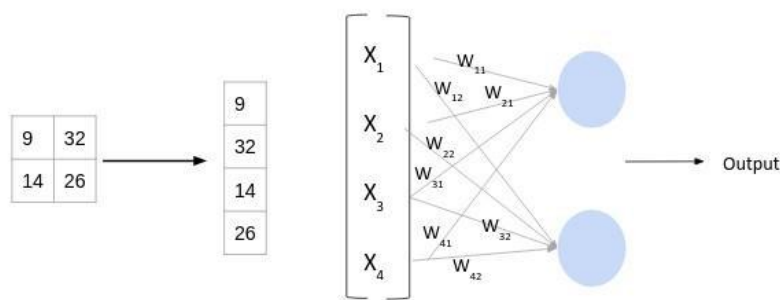


Once the data is converted into a 1D array, it is sent to the fully connected layer. All of these individual values are treated as separate features that represent the image. The fully connected layer performs two operations on the incoming data – **a linear transformation and a non-linear transformation.**
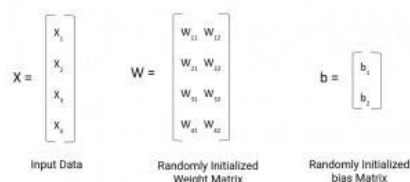
We first perform a linear transformation on this data. The equation for linear transformation is:

$$Z = W^T.X + b$$

Here, X is the input, W is weight, and b (called bias) is a constant. Note that the W in this case will be a matrix of (randomly initialized) numbers. Can you guess what would be the size of this matrix?



Considering the size of the matrix is (m, n) – m will be equal to the number of features or inputs for this layer. Since we have 4 features from the convolution layer, m here would be 4. The value of n will depend on the number of neurons in the layer. For instance, if we have two neurons, then the shape of weight matrix will be (4, 2):

Having defined the weight and bias matrix, let us put these in the equation for linear transformation:

$$Z = W^T . X + b$$

$$Z = \begin{bmatrix} W_{11} & W_{21} & W_{31} & W_{41} \\ W_{12} & W_{22} & W_{32} & W_{42} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$Z_{2x2} = \begin{bmatrix} W_{11}X_1 + W_{21}X_2 + W_{31}X_3 + W_{41}X_4 \\ W_{12}X_1 + W_{22}X_2 + W_{32}X_3 + W_{42}X_4 \end{bmatrix}$$

Now, there is one final step in the forward propagation process – the non-linear transformations. Let us understand the concept of non-linear transformation and it's role in the forward propagation process.

# Non-Linear transformation

The linear transformation alone can not capture complex relationships. Thus, we introduce an additional component in the network which adds non-linearity to the data. This new component in the architecture is called the **activation function**.

There are a number of activation functions that you can use – here is the complete list:

- [Fundamentals of Deep Learning – Activation Functions and When to Use Them?](#)

These activation functions are added at each layer in the neural network. The activation function to be used will depend on the type of problem you are solving.

We will be working on a binary classification problem and will use the Sigmoid activation function. Let's quickly look at the mathematical expression for this:

```
f(x) = 1/(1+e^-x)
```

The range of a Sigmoid function is between 0 and 1. This means that for any input value, the result would always be in the range (0, 1). A Sigmoid function is majorly used for binary classification problems and we will use this for both convolution and fully-connected layers.

Let's quickly summarize what we've covered so far.

# Forward Propagation Summary

**Step 1:** Load the input images in a variable (say X)

**Step 2:** Define (randomly initialize) a filter matrix. Images are convolved with the filter

```
Z1 = X * ƒ
```

Having defined the weight and bias matrix, let us put these in the equation for linear transformation :

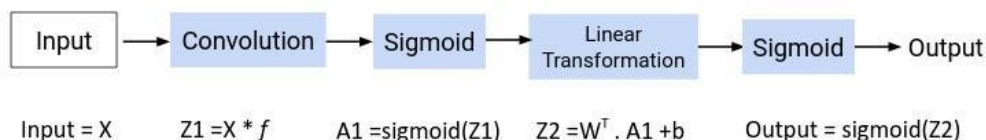**Step 3:** Apply the Sigmoid activation function on the result

```
A = sigmoid(Z1)
```

**Step 4:** Defin e (ran domly in itialize) weigh t an d bias matrix. A pply lin ear tran sformation on th e valu es

```
Z2 = Wᵀ.A + b
```

**Step 5:** A pply th e Sigmoid fu n ction on th e data. T h is will be th e fin al ou tpu t

```
O = sigmoid(Z2)
```
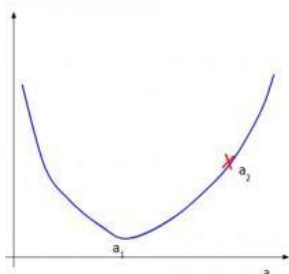


Now th e qu estion is – h ow are th e valu es in th e filter decided? T h e C NN model treats th ese valu es as parameters, wh ich are ran domly in itialized an d learn ed du rin g th e train in g process. W e will an swer th is in th e n ext section .

# Convolutional Neural Network (CNN): Backward Propagation

Du rin g th e forward propagation process, we ran domly in itialized th e weigh ts, biases an d filters. T h ese valu es are treated as parameters from th e con volu tion al n eu ral n etwork algorith m. **In the backwar d pr opagati on pr ocess, the model tr i es to update the par ameter s such that the ov er al l pr edi cti ons ar e mor e accur ate.**

F or u pdatin g th ese parameters, we u se th e g r a d i e n t  d e s c e n t  t e c h n i q u e . Let u s u n derstan d th e con cept of gradien t descen t with a simple example.
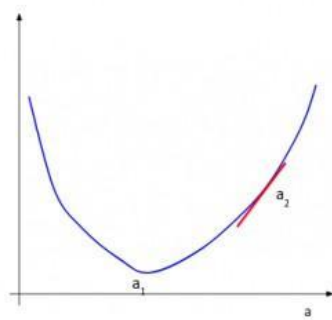
C on sider th at followin g in th e cu rve for ou r loss fu n ction wh ere we h ave a parameter a:



Du rin g th e ran dom in itialization of th e parameter, we get th e valu e of a as $a_2$ . I t is clear from th e pictu re th at th e min imu m valu e of loss is at $a_1$ an d n ot $a_2$ . T h e gradien t descen t tech n iqu e tries to fin d th is valu e of parameter (a) at wh ich th e loss is min imu m.

We understand that we need to update the value $a_2$ and bring it closer to $a_1$. To decide the direction of movement, i. e. whether to increase or decrease the value of the parameter, we calculate the gradient or

slope at the current point.

Based on the value of the gradient, we can determine the updated parameter values. When the slope is negative, the value of the parameter will be increased, and when the slope is positive, the value of the parameter should be decreased by a small amount.

Here is a generic equation for updating the parameter values:

```
new_parameter = old_parameter - (learning_rate * gradient_of_parameter)
```

The learning rate is a constant that controls the amount of change being made to the old value. The slope or the gradient determine the direction of the new values, that is, should the values be increased or decreased. So, we need to find the gradients, that is, change in error with respect to the parameters in order to update the parameter values.

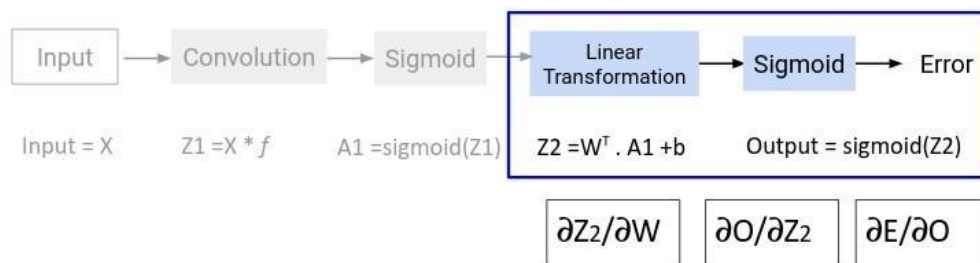*If you want to read about the gradient descent technique in detail, you can go through the below article:*

- [Introduction to Gradient Descent Algorithm](#)

We know that we have three parameters in a CNN model – weights, biases and filters. Let us calculate the gradients for these parameters one by one.

# Backward Propagation: Fully Connected Layer

As discussed previously, the fully connected layer has two parameters – weight matrix and bias matrix. Let us start by calculating the change in error with respect to weights – $\partial E / \partial W$.

Since the error is not directly dependent on the weight matrix, we will use the concept of chain rule to find this value. The computation graph shown below will help us define $\partial E / \partial W$ :

Backward Propagation (Fully Connected layer)

$$\partial E/\partial W = \partial E/\partial O \cdot \partial O/\partial Z_2 \cdot \partial z/\partial W$$

W e will fin d th e valu es of th ese derivatives separately.

## 1. Change in error with respect to output

Su ppose th e actu al valu es for th e data are den oted as y' an d th e predicted ou tpu t is represen ted as O. T h en th e error wou ld be given by th is equ ation :

$$E = (y' - O)^2/2$$

I f we differen tiate th e error with respect to th e ou tpu t, we will get th e followin g equ ation :

$$\partial E/\partial O = -(y'-O)$$

## 2. Change in output with respect to $Z_2$ (linear transformation output)

To fin d th e derivative of ou tpu t O with respect to $Z_2$ , we mu st first defin e O in terms of $Z_2$ . I f you look at th e compu tation graph from th e forward propagation section above, you wou ld see th at th e ou tpu t is simply th e sigmoid of $Z_2$ . T h u s, $\partial O/\partial Z_2$ is effectively th e derivative of Sigmoid. Recall th e equ ation for th e Sigmoid fu n ction :

$$f(x) = 1/(1+e^{-x})$$

T h e derivative of th is fu n ction comes ou t to be:

$$f'(x) = (1+e^{-x})^{-1}[1-(1+e^{-x})^{-1}] \quad f'(x) = sigmoid(x)(1-sigmoid(x)) \quad \partial O/\partial Z_2 = (O)(1-O)$$

You can read abou t th e complete derivation of th e Sigmoid fu n ction [here](#) .

# 3. Change in $Z_2$ with respect to W (Weights)

The value $Z_2$ is the result of the linear transformation process. Here is the equation of $Z_2$ in terms of weights:

```
Z2 = W^T.A1  + b
```

On differentiating $Z_2$ with respect to W, we will get the value $A_1$ itself:

```
∂Z2/∂W  = A1
```

Now that we have the individual derivations, we can use the chain rule to find the change in error with respect to weights:

```
∂E/∂W = ∂E/∂O . ∂O/∂Z2.  ∂Z2/∂W  ∂E/∂W = -(y'-o)  . sigmoid'.  A1
```

The shape of $\partial E/\partial W$ will be the same as the weight matrix W. We can update the values in the weight matrix using the following equation :
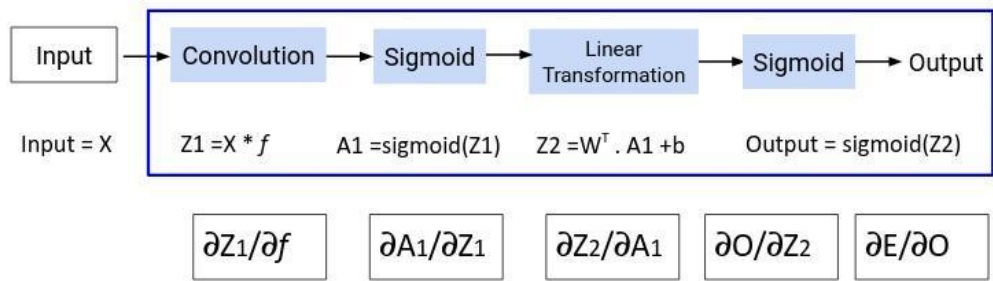
```
W_new = W_old - lr*∂E/∂W
```

Updating the bias matrix follows the same procedure. Try to solve that yourself and share the final equations in the comments section below!

# Backward Propagation: Convolution Layer

For the convolution layer, we had the filter matrix as our parameter. During the forward propagation process, we randomly initialized the filter matrix. We will now update these values using the following equation :

```
new_parameter = old_parameter - (learning_rate * gradient_of_parameter)
```

To update the filter matrix, we need to find the gradient of the parameter – dE/ df. Here is the computation graph for backward propagation :

Input → Convolution → Sigmoid → Linear Transformation → Sigmoid → Output

Input = X    Z1 = X * $f$    A1 = sigmoid(Z1)    Z2 = $W^T$ . A1 + b    Output = sigmoid(Z2)

| $\partial Z_1/\partial f$ | $\partial A_1/\partial Z_1$ | $\partial Z_2/\partial A_1$ | $\partial O/\partial Z_2$ | $\partial E/\partial O$ |

Backward Propagation (Convolution layer)

Input → Convolution → Sigmoid → Linear Transformation → Sigmoid → Output

Input = X    Z1 = X * $f$    A1 = sigmoid(Z1)    Z2 = $W^T$ . A1 + b    Output = sigmoid(Z2)

$\partial Z_1/\partial f$    $\partial A_1/\partial Z_1$    $\partial Z_2/\partial A_1$    $\partial O/\partial Z_2$    $\partial E/\partial O$

From the above graph we can define the derivative $\partial E/\partial f$ as:

```
∂E/∂f = ∂E/∂O.∂O/∂Z2.∂Z2/∂A1 .∂A1/∂Z1.∂Z1/∂f
```

We have already determined the values for $\partial E/\partial O$ and $\partial O/\partial Z2$. Let us find the values for the remaining derivatives.

## 1. Change in $Z_2$ with respect to $A_1$

To find the value for $\partial Z_2/\partial A_1$, we need to have the equation for $Z_2$ in terms of $A_1$:

```
Z2 = Wᵀ.A1  + b
```

On differentiating the above equation with respect to $A_1$, we get $W^T$ as the result:

```
∂Z2/∂A1  = Wᵀ
```

## 2. Change in $A_1$ with respect to $Z_1$

The next value that we need to determine is $\partial A_1/\partial Z_1$. Have a look at the equation of $A_1$

```
A1 = sigmoid(Z1)
```

This is simply the Sigmoid function. The derivative of Sigmoid would be:

```
∂A1/∂Z1  = (A1)(1-A1)
```

## 3. Change in $Z_1$ with respect to filter $f$

Finally, we need the value for $\partial Z_1/\partial f$. Here's the equation for $Z_1$

```
Z1 = X * f
```

Differentiating $Z$ with respect to $X$ will simply give us $X$:

```
∂Z1/∂f  = X
```

Now that we have all the required values, let's find the overall change in error with respect to the filter:

$$\partial E/\partial f \quad = \partial E/\partial O.\partial O/\partial Z_2.\partial Z_2/\partial A_1 .\partial A_1/\partial Z_1 \ * \ \partial Z_1/\partial f$$

Notice that in the equation above, the value $(\partial E/\partial O. \partial O/\partial Z_2. \partial Z_2/\partial A_1 . \partial A_1/\partial Z)$ is convolved with $\partial Z_1/\partial f$ instead of using a simple dot product. Why? The main reason is that during forward propagation, we perform a convolution operation for the images and filters.

This is repeated in the backward propagation process. Once we have the value for $\partial E/\partial f$, we will use this value to update the original filter value:

$$f = f - lr*(\partial E/\partial f)$$

This completes the back propagation section for convolutional neural networks. It's now time to code!

# CNN from Scratch using NumPy

Excited to get your hands dirty and design a convolutional neural network from scratch? The wait is over!

We will start by loading the required libraries and dataset. Here, **we will be using the MNIST dataset which is present within the *keras.datasets* library.**

For the purpose of this tutorial, we have selected only the first 200 images from the dataset. Here is the distribution of classes for the first 200 images:

```python
1   # importing  required  libraries
2   import  numpy  as np
3   import  pandas  as pd
4   from tqdm import  tqdm
5   from keras.datasets import  mnist
6
7   # loading  dataset
8   (x_train,  y_train),  (x_test,  y_test)  = mnist.load_data()
9
10  # selecting  a subset  of data (200 images)
11  x_train  = x_train[:200]
12  y = y_train[:200]
13
14  X = x_train.T
15  X = X/255
16
17  y.resize((200,1))
18  y = y.T
19
20  #checking  value
21  pd.Series(y[0]).value_counts()
```
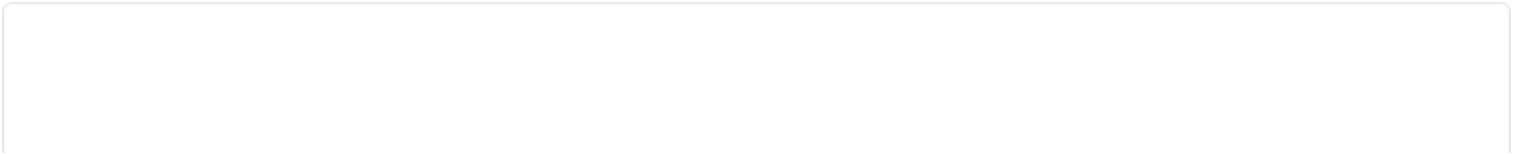
1 26 9 23 7 21 4 21 3 21 0 21 2 20 6 19 8 15 5 13 dtype:  int64

As you can see, we have ten classes here – 0 to 9. **This is a multi-class classification problem**. For now, we will start with building a simple CNN model for a binary classification problem:

```python
1   # converting  into binary  classification
2   for i in range(y.shape[1]):
```

```
3        if y[0][i]  >4:
4            y[0][i]  = 1
```

```
5        else:
6            y[0][i]  = 0
7
8    #checking  value  counts
9    pd.Series(y[0]).value_counts()
```

```
0 109 1 91 dtype:  int64
```

We will now initialize the filters for the convolution operation :

```
1    # initializing filter
2    f=np.random.uniform(size=(3,5,5))
3    f = f.T
4
5    print('Filter 1', '\n',  f[:,:,0],  '\n')
6    print('Filter 2', '\n',  f[:,:,1],  '\n')
7    print('Filter 3', '\n',  f[:,:,2],  '\n')
```

Let us quickly check the shape of the loaded images, target variable and the filter matrix:

```
X.shape,  y.shape,  f.shape
```

```
((28, 28, 200), (1, 200), (5, 5, 3))
```

We have 200 images of dimensions (28, 28) each . For each of these 200 images, we have one class specified and hence the shape of y is (1, 200). Finally, we defined 3 filters, each of dimensions (5, 5).

The next step is to prepare the data for the convolution operation . As we discussed in the forward propagation section , the image-filter convolution can be treated as a single image being divided into multiple patches:
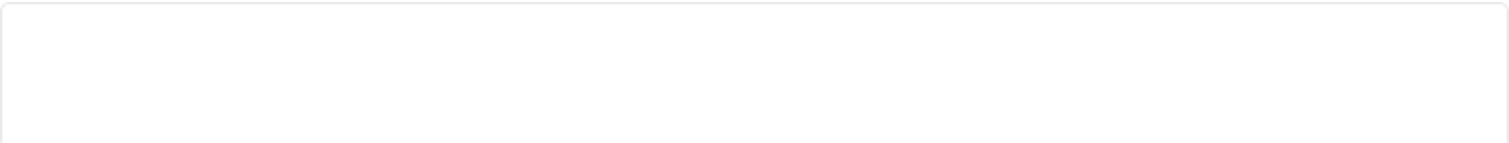


For every single image in the data, we will create smaller patches of the same dimension as the filter matrix, which is (5, 5). Here is the code to perform this task :

```
1    # Generating  patches  from images
2    new_image  = []
```

```
3
4    # for number  of images
```

```
5    for k in range(X.shape[2]):
6        # sliding  in horizontal  direction
7        for i in range(X.shape[0]-f.shape[0]+1):
8            # sliding  in vertical  direction
9            for j in range(X.shape[1]-f.shape[1]+1):
10               new_image.append(X[:,:,k][i:i+f.shape[0],j:j+f.shape[1]])
11
12   # resizing  the generated  patches  as per number  of images
13   new_image  = np.array(new_image)
14   new_image.resize((X.shape[2],int(new_image.shape[0]/X.shape[2]),new_image.shape[1],new_image.shape[2]))
15   new_image.shape
```

(200, 576, 5, 5)

We h ave everyth in g we n eed  for th e forward propagation process of th e con volu tion  layer. Movin g on  to th e n ext  section for  forward propagation ,  we  n eed  to  in itialize th e  weigh t matrix  for  th e  fu lly con n ected layer. T h e size  of th e weigh t matrix will be  $(m,\ n\ )$ – wh ere  m is  th e n u mber of featu res as in pu t an d n will be  th e n u mber of n eu ron s in th e  layer.

W h at wou ld  be  th e  n u mber of  featu res? W e  k n ow  th at we  can determin e  th e  sh ape of  th e  ou tpu t  image u sin g th is  formu la:

$$((n-f+1)\ \ ,\ (n-f+1))$$

Sin ce  th e  fu lly con n ected  layer on ly tak es 1D in pu t,  we  will flatten  th is  matrix.  T h is  mean s th e  n u mber of featu res or  in pu t valu es will be:

$$((n-f+1)\ \ x\ (n-f+1)\ \ x\ num\_of\_filter)$$

Let u s in itialize th e weigh t matrix:

```
1    # number  of features  in data set
2    s_row  = X.shape[0]  - f.shape[0]  + 1
3    s_col  = X.shape[1]  - f.shape[1]  + 1
4    num_filter  = f.shape[2]
5
6    inputlayer_neurons = (s_row)*(s_col)*(num_filter)
7    output_neurons = 1
8
9    # initializing weight
10   wo=np.random.uniform(size=(inputlayer_neurons,output_neurons))
```

F in ally,  we  will write  th e code for  th e activation  fu n ction  wh ich  we  will be  u sin g for  th e con volu tion  n eu ral n etwork  arch itectu re.

Now, we h ave already con ver ted  th e origin al problem in to a bin ary classification problem. Hen ce, we will be u sin g  Sigmoid  as  ou r  activation   fu n ction .  W e  h ave  already   discu ssed  th e  math ematical   equ ation   for Sigmoid an d its derivative. Here  is th e pyth on  code for th e same:

```
1    # defining  the Sigmoid  Function
2    def sigmoid  (x):
3        return  1/(1 + np.exp(-x))
4
5    # derivative  of Sigmoid  Function
6    def derivatives_sigmoid(x):
7        return  x * (1 - x)
```

Great! We have all the elements we need for the forward propagation process. Let us put these code blocks together.

First, we perform the convolution operation on the patches created. After the convolution, the results are stored in the form of a list, which is converted into an array of dimension (200, 3, 576). Here

- 200 is the number of images
- 576 is the number of patches created, and
- 3 is the number of filters we used

After the convolution operation, we apply the Sigmoid activation function :

```python
# generating output of convolution layer
filter_output = []
# for each image
for i in range(len(new_image)):
    # apply each filter
    for k in range(f.shape[2]):
        # do element wise multiplication
        for j in range(new_image.shape[1]):
            filter_output.append((new_image[i][j]*f[:,:,k]).sum())

filter_output = np.resize(np.array(filter_output), (len(new_image),f.shape[2],new_image.shape[1]))

# applying activation over convolution output
filter_output_sigmoid = sigmoid(filter_output)

filter_output.shape, filter_output_sigmoid.shape
```

((200, 3, 576), (200, 3, 576))

After the convolution layer, we have the fully connected layer. We know that the fully connected layer will only have 1D inputs. So, we first flatten the results from the previous layer using the *reshape* function.
Then, we apply the linear transformation and activation function on this data:

```python
# generating input for fully connected layer
filter_output_sigmoid = filter_output_sigmoid.reshape((filter_output_sigmoid.shape[0],filter_output_sigmoid.shape[1]*filter_outp

filter_output_sigmoid = filter_output_sigmoid.T

# Linear trasnformation for fully Connected Layer
output_layer_input= np.dot(wo.T,filter_output_sigmoid)
output_layer_input = (output_layer_input - np.average(output_layer_input))/np.std(output_layer_input)

# activation function
output  = sigmoid(output_layer_input)
```

It's time to start the code for backward propagation. Let's define the individual derivatives for backward propagation of the fully connected layer. Here are the equations we need:

$$E = (y' - O)^2/2 \quad \partial E/\partial O = -(y'-O) \quad \partial O/\partial Z2 = (O)(1-O) \quad \partial Z2/\partial W = A1$$

Let us code this in Python :

```python
#Error
error = np.square(y-output)/2

```

```
4    #Error  w.r.t  Output  (Gradient)
5    error_wrt_output = -(y-output)
```

We have the individual derivatives from the previous code block. We can now find the overall change in error w.r.t. weight using the chain rule. Finally, we will use this gradient value to update the original weight matrix:

```
W_new = W_old - lr*∂E/∂W
```

So far we have covered back propagation for the fully connected layer. This covers updating the weight matrix. Next, we will look at the derivatives for back propagation for the convolutional layer and update the filters:

$$\partial E/\partial f = \partial E/\partial O.\partial O/\partial Z_2.\partial Z_2/\partial A_1 .\partial A_1/\partial Z_1 * \partial Z_1/\partial f \quad \partial E/\partial O = -(y'-O) \quad \partial O/\partial Z_2 = (O)(1-O) \quad \partial Z_2/\partial A_1 = W^T \quad \partial A_1/\partial Z_1 = A_1(1-A_1)$$
$$\partial Z_1/\partial f = X$$

We will code the first four equations in Python and calculate the derivative using the *np.dot* function. Post that we need to perform a convolution operation using $\partial Z_1/\partial f$:

We now have the gradient value. Let us use it to update the error:

```
1   for i in range(f.shape[2]):
2       f[:,:,i]  = f[:,:,i]  - lr*filter_update_array[i]
```

# End Notes

Take a deep breath – that was a lot of learning in one tutorial! Convolutional neural networks can appear to be slightly complex when you're starting out but once you get the hang of how they work, you'll feel ultra confident in yourself.

I had a lot of fun writing about what goes on under the hood of these CNN models we see everywhere these days. We covered the mathematics behind these CNN models and learned how to implement them from scratch, using just the NumPy library.

I want you to explore other concepts related to CNN, such as the padding and pooling techniques. Implement them in the code as well and share your ideas in the comments section below.

You should also try your hand at the below hackathons to practice what you've learned:

- [Practice Problem: Identify the Apparels](#)
- [Practice Problem: Identify the Digits](#)