# UE21CS343BB2

# Topics in Deep Learning

**Dr. Shylaja S S**
Director of Cloud Computing & Big Data (CCBD), Centre for Data Sciences & Applied Machine Learning (CDSAML)
Department of Computer Science and Engineering
**shylaja.sharath@pes.edu**

**Ack: Divya K,
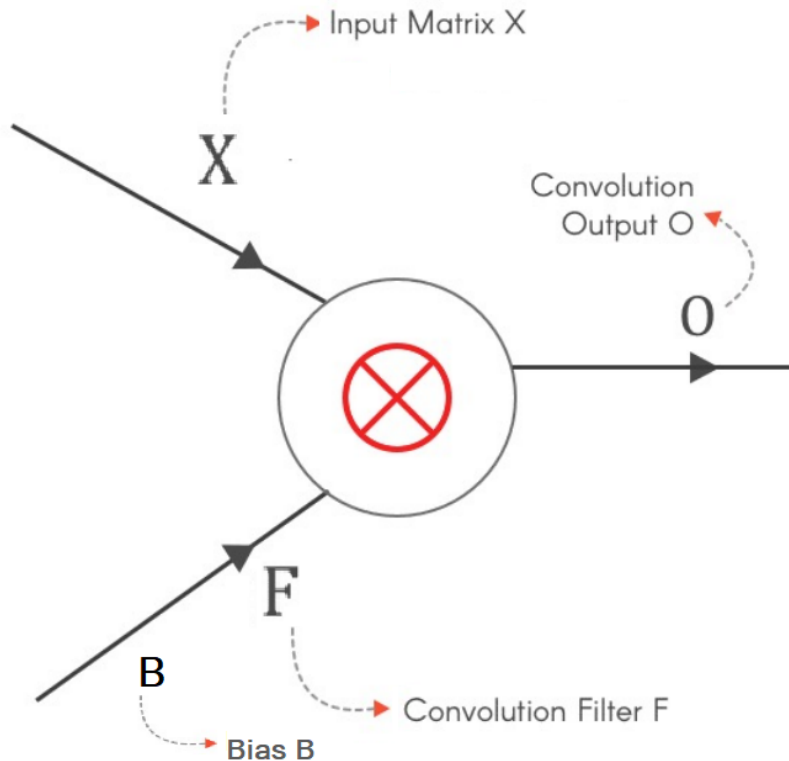Teaching Assistant**

**Overview of lecture**

- Recap: Backpropagation
- Convolution Forward Pass
- Convolution Backward Pass
  - ➢ Calculation of loss gradient w.r.t filter
  - ➢ Calculation of loss gradient w.r.t bias
  - ➢ Calculation of loss gradient w.r.t input

- Conclusion: Backpropagation in Convolution Layer

**Recap: Backpropagation**

➢ Backpropagation is an algorithm used to train neural networks by adjusting the weights of the network based on the error between the predicted output and the actual output.

➢ Backpropagation calculates the gradient of the loss function with respect to each parameter in the network and updates the network parameters in such a way that it minimizes the loss function.

➢ In CNNs the loss gradient is computed w.r.t the input and also w.r.t the filter, w.r.t the bias.

Convolution between Input X and Filter F, gives us an output O. This can be represented as:



$$\begin{bmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{bmatrix} = \text{Convolution} \left( \begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}, \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \right) + \text{Bias B}$$

Output O          Input X          Filter F
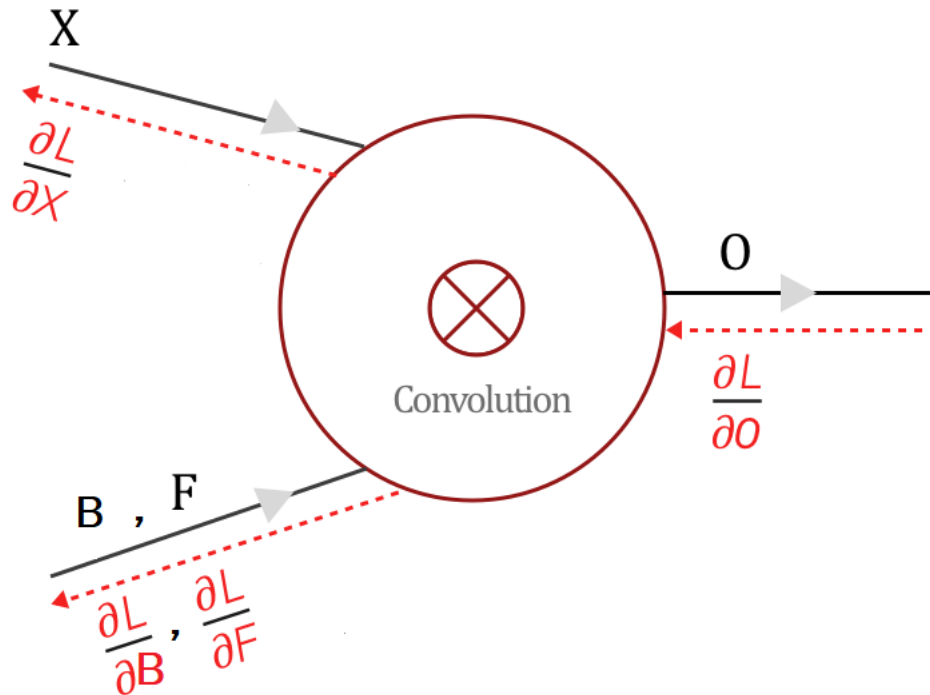
Applying the convolution operation, we get:



$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22} + B$$

$$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22} + B$$

$$O_{21} = X_{21}F_{11} + X_{22}F_{12} + X_{31}F_{21} + X_{32}F_{22} + B$$

$$O_{22} = X_{22}F_{11} + X_{23}F_{12} + X_{32}F_{21} + X_{33}F_{22} + B$$

## Convolution Backward Pass



In ANNs, we update the weights as
W = W - α*∂L/∂W
In CNNs, we update the network parameters as:
F = F - α*∂**L**/∂**F**
B = B - α*∂**L**/∂**B**
(where α is the learning parameter)
Since X is the output of the previous layer, ∂**L**/∂**X** becomes the loss gradient for the previous layer. So, we need to calculate ∂**L**/∂**F**, ∂**L**/∂**B** and ∂**L**/∂**X**.

➢ Calculation of loss gradient w.r.t the filter, $\partial L/\partial F$:

We can use the chain rule to obtain the gradient w.r.t the filter as shown in the equation.

$$\frac{\partial L}{\partial F} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial F}$$

Gradient to update Filter F

Loss Gradient from previous layer

Local Gradients

*For every element of F*

$$\frac{\partial L}{\partial F_i} = \sum_{k=1}^{M} \frac{\partial L}{\partial O_k} * \frac{\partial O_k}{\partial F_i}$$

**Convolution Backward Pass**

➢ Calculation of loss gradient w.r.t the filter, $\partial L/\partial F$:

On expanding the chain rule summation, we get:

$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{11}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{11}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{11}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{11}}$$

$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{12}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{12}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{12}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{12}}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{21}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{21}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{21}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{21}}$$

$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{22}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{22}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{22}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{22}}$$

➢ Calculation of loss gradient w.r.t the filter, $\partial$**L**/$\partial$**F:**

To calculate $\partial$**O**/$\partial$**F**:

From these equations,

$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22} + B$

$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22} + B$

$O_{21} = X_{21}F_{11} + X_{22}F_{12} + X_{31}F_{21} + X_{32}F_{22} + B$

$O_{22} = X_{22}F_{11} + X_{23}F_{12} + X_{32}F_{21} + X_{33}F_{22} + B$

$\longrightarrow$

We get:

$$\frac{\partial O_{11}}{\partial F_{11}} = X_{11}, \frac{\partial O_{11}}{\partial F_{12}} = X_{12}, \frac{\partial O_{11}}{\partial F_{21}} = X_{21}, \frac{\partial O_{11}}{\partial F_{22}}$$

$$\frac{\partial O_{12}}{\partial F_{11}} = X_{12}, \frac{\partial O_{12}}{\partial F_{12}} = X_{13}, \frac{\partial O_{12}}{\partial F_{21}} = X_{22}, \frac{\partial O_{1}}{\partial F_{2}}$$

$$\frac{\partial O_{21}}{\partial F_{11}} = X_{21}, \frac{\partial O_{21}}{\partial F_{12}} = X_{22}, \frac{\partial O_{21}}{\partial F_{21}} = X_{31}, \frac{\partial O_{2}}{\partial F_{2}}$$

$$\frac{\partial O_{22}}{\partial F} = X_{22}, \frac{\partial O_{22}}{\partial F} = X_{23}, \frac{\partial O_{22}}{\partial F} = X_{32}, \frac{\partial O_{2}}{\partial F}$$

**Convolution Backward Pass**

➢ Calculation of loss gradient w.r.t the filter, $\partial L/\partial F$:

On substituting the values of the local gradient, we get:

$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}} * X_{11} + \frac{\partial L}{\partial O_{12}} * X_{12} + \frac{\partial L}{\partial O_{21}} * X_{21} + \frac{\partial L}{\partial O_{22}} * X_{22}$$

$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}} * X_{12} + \frac{\partial L}{\partial O_{12}} * X_{13} + \frac{\partial L}{\partial O_{21}} * X_{22} + \frac{\partial L}{\partial O_{22}} * X_{23}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}} * X_{21} + \frac{\partial L}{\partial O_{12}} * X_{22} + \frac{\partial L}{\partial O_{21}} * X_{31} + \frac{\partial L}{\partial O_{22}} * X_{32}$$

$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}} * X_{22} + \frac{\partial L}{\partial O_{12}} * X_{23} + \frac{\partial L}{\partial O_{21}} * X_{32} + \frac{\partial L}{\partial O_{22}} * X_{33}$$

## Convolution Backward Pass

➢ Calculation of loss gradient w.r.t the filter, $\partial L/\partial F$:

If we look closely, this can be represented as a **convolution operation between input X** and **loss gradient $\partial L/\partial O$** as shown below:

$$
\begin{bmatrix} \frac{\partial L}{\partial F_{11}} & \frac{\partial L}{\partial F_{12}} \\ \frac{\partial L}{\partial F_{21}} & \frac{\partial L}{\partial F_{22}} \end{bmatrix} = \text{Convolution} \left( \begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}, \begin{bmatrix} \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \end{bmatrix} \right)
$$

where

$$
\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix} = \text{Input X}
$$

$$
\begin{bmatrix} \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \end{bmatrix} = \frac{\partial L}{\partial O} \quad \text{Loss gradient from previous layer}
$$

$\partial L/\partial F$ = Convolution of input matrix X and loss gradient $\partial L/\partial O$

$$
\frac{\partial L}{\partial F} = \text{conv}(X, \frac{\partial L}{\partial O})
$$

➢ Calculation of loss gradient w.r.t the bias, $\partial L/\partial B$:

We can use the chain rule to obtain the gradient w.r.t the bias as shown in the equation.

$$\frac{\partial L}{\partial B} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial B}$$

Gradient to update **Bias B**

Loss Gradient from previous layer

Local Gradients

$$\frac{\partial L}{\partial B} = \sum_{k=1}^{M} \frac{\partial L}{\partial O_k} * \frac{\partial O_k}{\partial B}$$

➢ Calculation of loss gradient w.r.t the bias, $\partial \textbf{L}/\partial \textbf{B}$:

To calculate $\partial \textbf{O}/\partial \textbf{B}$, we just partially derive $\textbf{O}_{11}$, $\textbf{O}_{12}$, $\textbf{O}_{21}$, and $\textbf{O}_{22}$ with respect to $\textbf{B}$. Since there is only one $\textbf{B}$ term in each $\textbf{O}$ term (as shown), the partial differentiation just returns 1.

$$\frac{\partial O_{11}}{\partial B} = 1$$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22} + B$$

$$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22} + B$$

$$\frac{\partial O_{12}}{\partial B} = 1$$

$$O_{21} = X_{21}F_{11} + X_{22}F_{12} + X_{31}F_{21} + X_{32}F_{22} + B$$

$$\frac{\partial O_{21}}{\partial B} = 1$$

$$O_{22} = X_{22}F_{11} + X_{23}F_{12} + X_{32}F_{21} + X_{33}F_{22} + B$$

$$\frac{\partial O_{22}}{\partial B} = 1$$

**Convolution Backward Pass**

➢ Calculation of loss gradient w.r.t the bias, $\partial L/\partial B$:

So $\partial L/\partial B$ is just equal to the summation of $\partial L/\partial O$ terms.

$$\frac{\partial L}{\partial B} = \sum \frac{\partial L}{\partial O_k} * \frac{^1 \partial O_k}{\partial B}$$

$$\frac{\partial L}{\partial B} = \frac{\partial L}{\partial O_{11}} + \frac{\partial L}{\partial O_{12}} + \frac{\partial L}{\partial O_{21}} + \frac{\partial L}{\partial O_{22}}$$

$$\frac{\partial L}{\partial B} = \sum_{k=1}^{M} \frac{\partial L}{\partial O_k}$$

$$\frac{\partial L}{\partial B} = \text{sum}(\frac{\partial L}{\partial O})$$

## Convolution Backward Pass

➢ Calculation of loss gradient w.r.t the input, $\partial L/\partial X$:

We can use the chain rule to obtain the gradient w.r.t the input as shown in the equation.

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial X}$$

Gradient to update input **X**

Loss Gradient from previous layer

Local Gradients

For every element of $X_i$

$$\frac{\partial L}{\partial X_i} = \sum_{k=1}^{M} \frac{\partial L}{\partial O_k} * \frac{\partial O_k}{\partial X_i}$$

**Convolution Backward Pass**

➢ Calculation of loss gradient w.r.t the input, $\partial L/\partial X$:

On expanding the chain rule summation and substituting the values of the local gradients, we get:

$$\frac{\partial L}{\partial X_{11}} = \frac{\partial L}{\partial O_{11}} * F_{11}$$

$$\frac{\partial L}{\partial X_{12}} = \frac{\partial L}{\partial O_{11}} * F_{12} + \frac{\partial L}{\partial O_{12}} * F_{11}$$

$$\frac{\partial L}{\partial X_{13}} = \frac{\partial L}{\partial O_{12}} * F_{12}$$

$$\frac{\partial L}{\partial X_{21}} = \frac{\partial L}{\partial O_{11}} * F_{21} + \frac{\partial L}{\partial O_{21}} * F_{11}$$

$$\frac{\partial L}{\partial X_{22}} = \frac{\partial L}{\partial O_{11}} * F_{22} + \frac{\partial L}{\partial O_{12}} * F_{21} + \frac{\partial L}{\partial O_{21}} * F_{12} + \frac{\partial L}{\partial O_{22}} * F_{11}$$

$$\frac{\partial L}{\partial X_{23}} = \frac{\partial L}{\partial O_{12}} * F_{22} + \frac{\partial L}{\partial O_{22}} * F_{12}$$

$$\frac{\partial L}{\partial X_{31}} = \frac{\partial L}{\partial O_{21}} * F_{21}$$
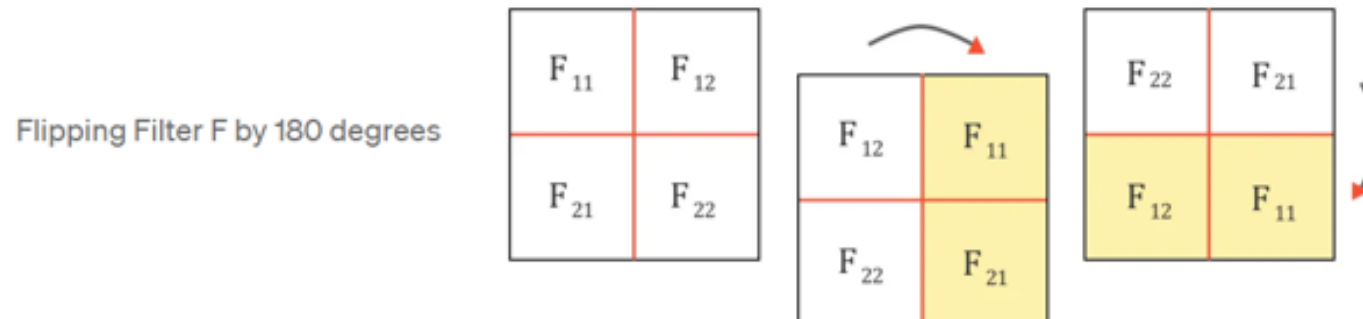
$$\frac{\partial L}{\partial X_{32}} = \frac{\partial L}{\partial O_{21}} * F_{22} + \frac{\partial L}{\partial O_{22}} * F_{21}$$

$$\frac{\partial L}{\partial X_{33}} = \frac{\partial L}{\partial O_{22}} * F_{22}$$

➤ Calculation of loss gradient w.r.t the input, $\partial L/\partial X$:

If we look closely, this can be represented as a **"full convolution" operation between flipped / 180° rotated Filter F and loss gradient $\partial L/\partial O$** as shown below:

❑ the term "full convolution" is often used interchangeably with "convolution with zero-padding." In the context of convolutional neural networks (CNNs), "full convolution" typically means performing a convolution operation with zero-padding applied to the input. Here, we mean padded loss gradient $\partial L/\partial O$.



Flipping Filter F by 180 degrees

## Convolution Backward Pass

➢ Calculation of loss gradient w.r.t the input, $\partial L/\partial X$:

Now, let us do a 'full' convolution between this flipped Filter F and $\partial L/\partial O$, as visualized below:



$$\frac{\partial L}{\partial X_{11}} = F_{11} * \frac{\partial L}{\partial O_{11}}$$

Filter F

Loss Gradient $\frac{\partial L}{\partial O}$

@pavisj

## Convolution Backward Pass

➢ Calculation of loss gradient w.r.t the input, *∂L/∂X:*

The full convolution above generates the values of *∂L/∂X* and hence we can represent *∂L/∂X* as:



$$\begin{bmatrix} \dfrac{\partial L}{\partial X_{11}} & \dfrac{\partial L}{\partial X_{12}} & \dfrac{\partial L}{\partial X_{13}} \\ \dfrac{\partial L}{\partial X_{21}} & \dfrac{\partial L}{\partial X_{22}} & \dfrac{\partial L}{\partial X_{23}} \\ \dfrac{\partial L}{\partial X_{31}} & \dfrac{\partial L}{\partial X_{32}} & \dfrac{\partial L}{\partial X_{33}} \end{bmatrix}$$

$\dfrac{\partial L}{\partial X}$

= Full Convolution

$$\left( \begin{bmatrix} F_{22} & F_{21} \\ F_{12} & F_{11} \end{bmatrix} , \begin{bmatrix} \dfrac{\partial L}{\partial O_{11}} & \dfrac{\partial L}{\partial O_{12}} \\ \dfrac{\partial L}{\partial O_{21}} & \dfrac{\partial L}{\partial O_{22}} \end{bmatrix} \right)$$

180-degree rotated Filter F        Loss Gradient $\dfrac{\partial L}{\partial O}$

*∂L/∂X can be represented as 'full' convolution between a 180-degree rotated Filter F and loss gradient ∂L/∂O*

$$\frac{\partial L}{\partial X} = \text{full-conv}(180° \text{ flipped F}, \frac{\partial L}{\partial O})$$

**OR**

$$\frac{\partial L}{\partial X} = \text{conv}(180° \text{ flipped F}, \text{padded}(\frac{\partial L}{\partial O}))$$

## Conclusion: Backpropagation in Convolution Layer

$$\frac{\partial L}{\partial F} = conv(X, \frac{\partial L}{\partial O})$$

$$\frac{\partial L}{\partial B} = sum(\frac{\partial L}{\partial O})$$

$$\frac{\partial L}{\partial X} = \text{full-conv}(180° \text{ flipped } F, \frac{\partial L}{\partial O})$$

**OR**

$$\frac{\partial L}{\partial X} = conv(180° \text{ flipped } F, padded(\frac{\partial L}{\partial O}))$$

In CNNs, we update the network parameters as:
F = F - α*∂L/∂F
B = B - α*∂L/∂B
(where α is the learning parameter)
Since X is the output of the previous layer, ∂L/∂X
becomes the loss gradient for the previous layer.

Note : This is done for a single filter F and stride = 1, CNNs have a lot more filters. The CNN Backpropagation operation with stride > 1 is identical to a stride = 1

## Acknowledgements

# UE21CS343BB2

## Topics in Deep Learning

**Dr. Shylaja S S**
Director of Cloud Computing & Big Data (CCBD), Centre for Data Sciences & Applied Machine Learning (CDSAML)
Department of Computer Science and Engineering
**shylaja.sharath@pes.edu**

**Ack: Divya K,
Teaching Assistant**