

Backpropagation

Convolutional Neural Networks

EE5179: Deep Learning for Imaging

Slides based on Pavithra Solai's article at <https://pavisj.medium.com/convolutions-and-backpropagations-46026a8f5d2c>,
Jefkine's article at <http://www.jefkine.com/general/2016/09/05/backpropagation-in-convolutional-neural-networks/>,
And many slides are from Rukayat Sadiq's slides on CNN backpropagation
https://deeplearning.cs.cmu.edu/F21/document/recitation/Recitation5/CNN_Backprop_Recitation_5_F21.pdf

Back-propagation algorithm for NN

1. Perform a feedforward pass
 - Computing activations L_1 , L_2 and so on ...

2. For each output unit i in layer L_4 (output layer), set

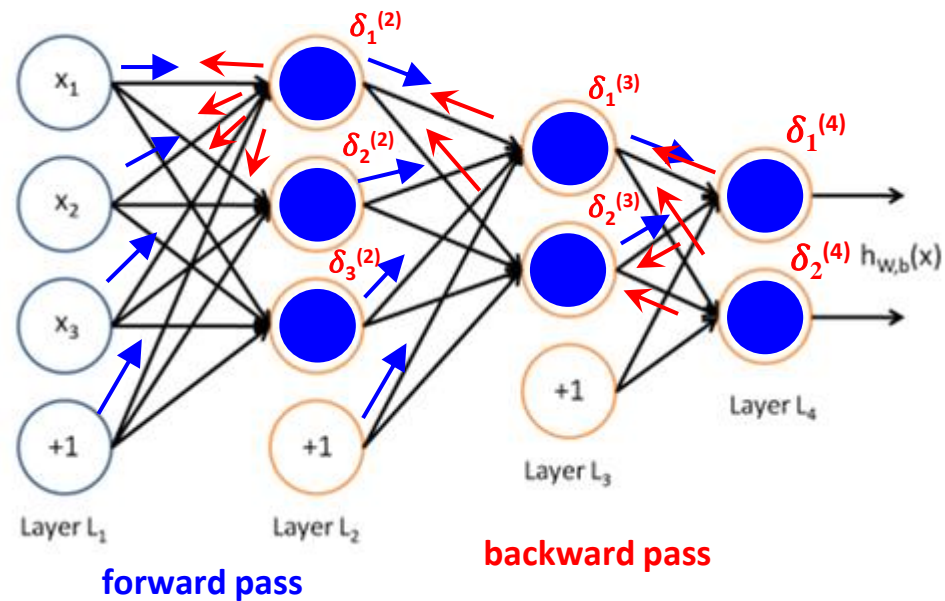
$$\delta^{(n_i)} = -(y - a^{(n_i)}) \bullet f'(z^{(n)})$$

3. Starting from last but one layer to 2nd layer;
 $l = n_l - 1, n_l - 2, \dots, 2$

$$\delta^{(l)} = ((W^{(l+1)})^T \delta^{(l+1)}) \bullet f'(z^{(l)})$$

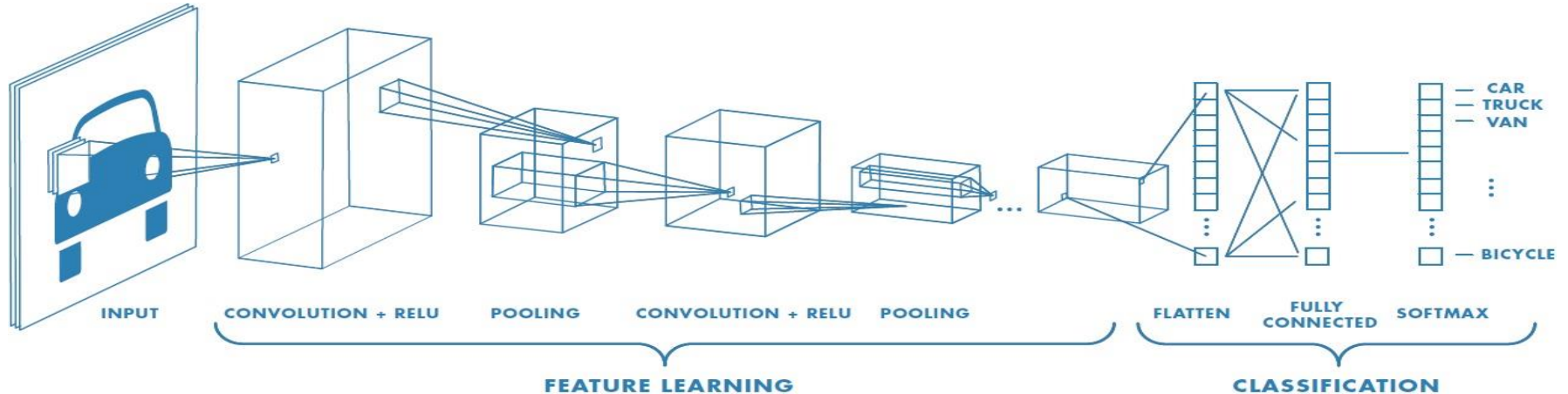
4. Compute the desired partial derivatives, as:

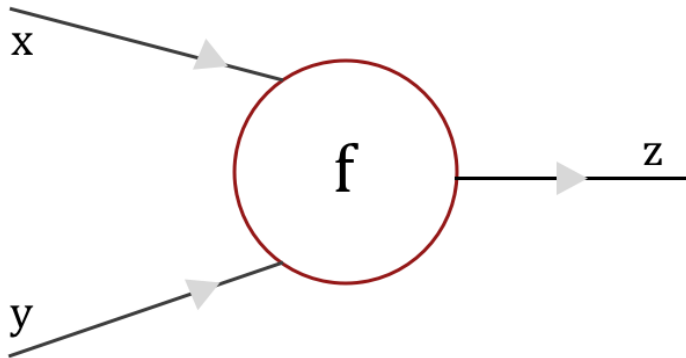
$$\begin{aligned}\nabla_{W^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} (a^{(l)})^T, \\ \nabla_{b^{(l)}} J(W, b; x, y) &= \delta^{(l+1)}.\end{aligned}$$



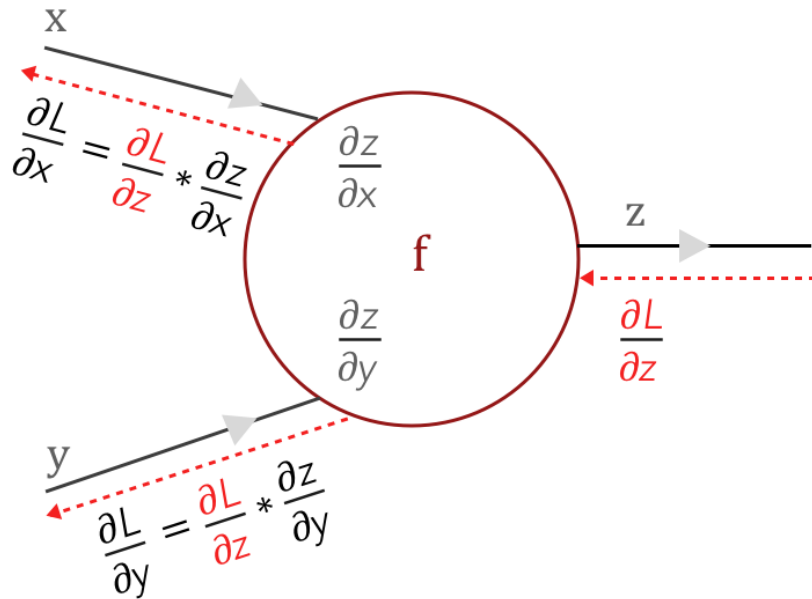
Review of Conv Nets

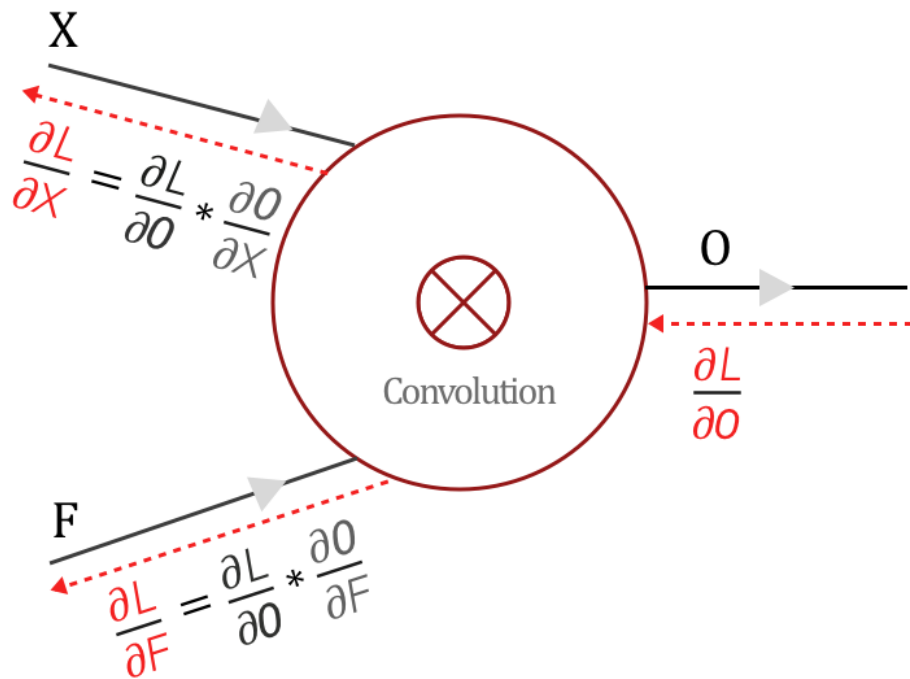
- Steps involved
 - Perform convolutions
 - Apply non-linearity
 - Pooling





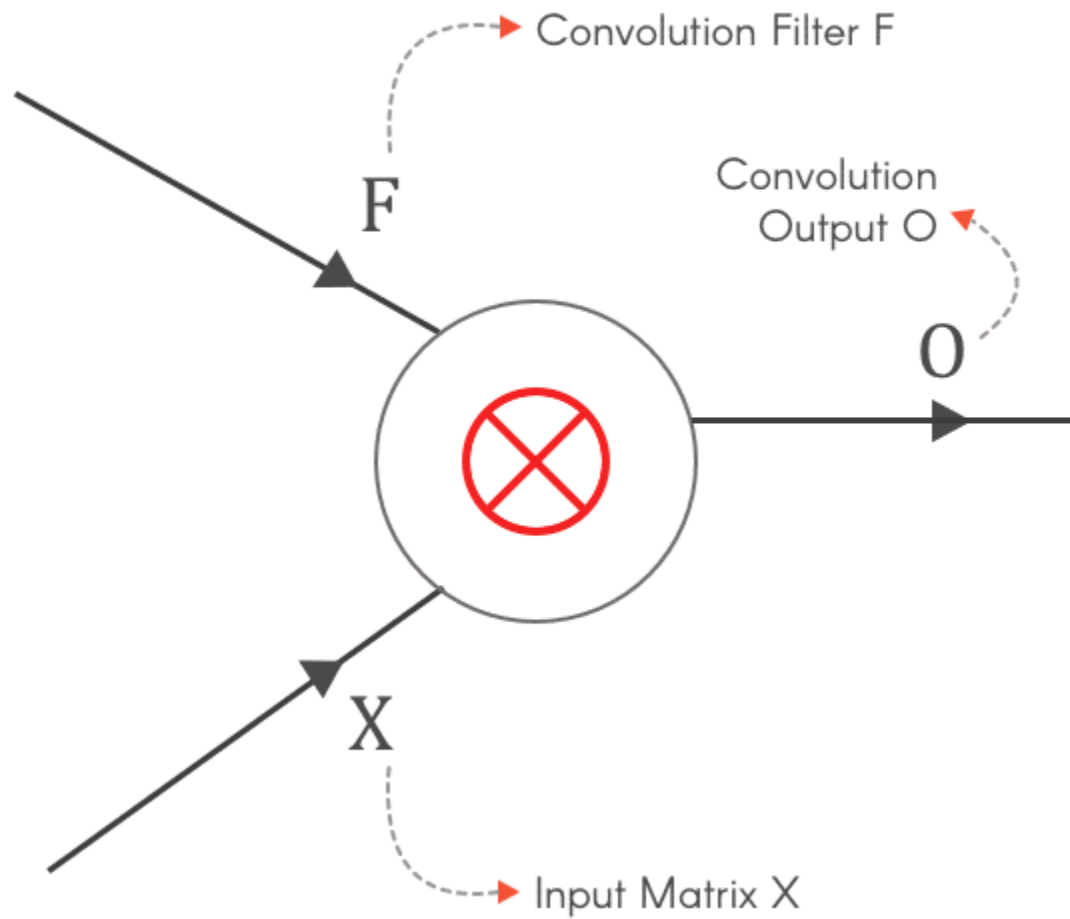
A simple function f which takes x and y as inputs and outputs z





$\frac{\partial O}{\partial X}$ & $\frac{\partial O}{\partial F}$ are local gradients

$\frac{\partial L}{\partial Z}$ is the loss from the previous layer which has to be backpropagated to other layers



Backpropagation in CNNs

- In the backward pass, we get the loss gradient with respect to the next layer
- In CNNs the loss gradient is computed w.r.t the input and also w.r.t the filter.

Convolution Backprop with single Stride

- To understand the computation of loss gradient w.r.t input, let us use the following example:

- Horizontal and vertical stride = 1

X_{11}	X_{12}	X_{13}
X_{21}	X_{22}	X_{23}
X_{31}	X_{32}	X_{33}

Input **X**

F_{11}	F_{12}
F_{21}	F_{22}

Filter **F**

Convolution Forward Pass

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

$$\begin{array}{|c|c|} \hline O_{11} & O_{12} \\ \hline O_{21} & O_{22} \\ \hline \end{array} = \text{Convolution} \left(\begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21} & X_{22} & X_{23} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array}, \begin{array}{|c|c|} \hline F_{11} & F_{12} \\ \hline F_{21} & F_{22} \\ \hline \end{array} \right)$$

Output O Input X Filter F

Convolution Forward Pass

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

X_{11}	X_{12}	X_{13}
X_{21}	X_{22}	X_{23}
X_{31}	X_{32}	X_{33}

Input X



F_{11}	F_{12}
F_{21}	F_{22}

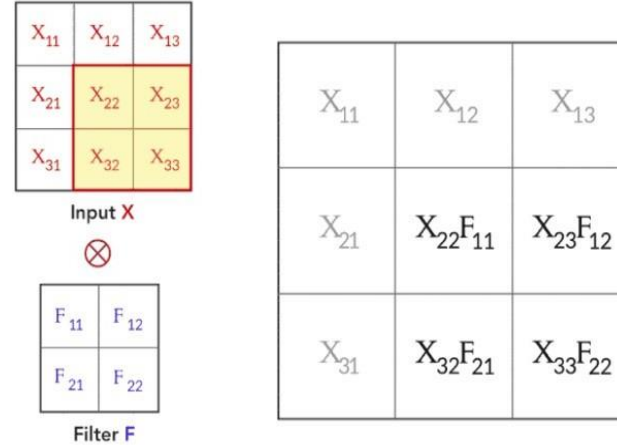
Filter F

$X_{11}F_{11}$	$X_{12}F_{12}$	X_{13}
$X_{21}F_{21}$	$X_{22}F_{22}$	X_{23}
X_{31}	X_{32}	X_{33}

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

Convolution Forward Pass

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



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

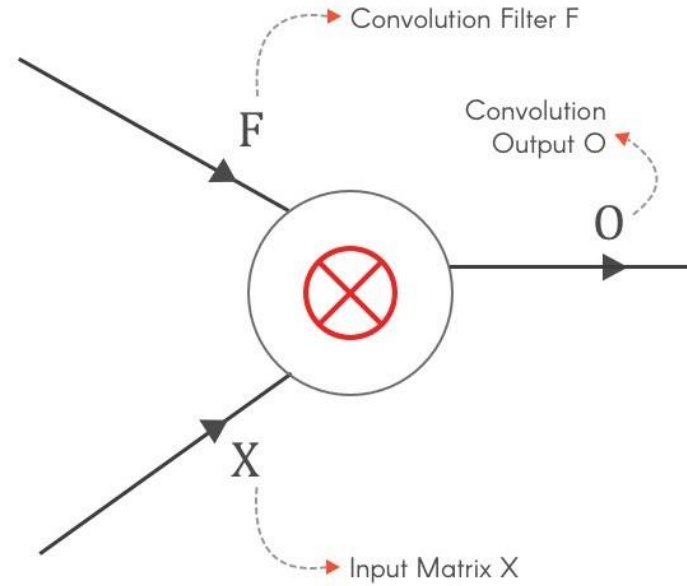
$$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22}$$

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

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

Loss gradient

- We want to calculate the gradients wrt to input 'X' and filter 'F'



Loss gradient w.r.t the filter

We can use the chain rule to obtain the gradient wrt the filter as shown in the equation.

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

Diagram illustrating the chain rule for the gradient of the loss with respect to the filter F :

- $\frac{\partial L}{\partial F}$ (red) is labeled "Gradient to update Filter F".
- $\frac{\partial L}{\partial O}$ is labeled "Loss Gradient from previous layer".
- $\frac{\partial O}{\partial F}$ is labeled "Local Gradients".

For every element of F

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

Loss gradient w.r.t the filter

We can expand the chain rule summation as:

For every element of F

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

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

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

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

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

Loss gradient w.r.t the filter

- Replacing the local gradients of the filter i.e, $\frac{\partial L}{\partial F_i}$ we get this:

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

where

$$\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix} = \text{Input X} \quad \begin{bmatrix} \frac{\partial L}{\partial \theta_{11}} & \frac{\partial L}{\partial \theta_{12}} \\ \frac{\partial L}{\partial \theta_{21}} & \frac{\partial L}{\partial \theta_{22}} \end{bmatrix} = \frac{\partial L}{\partial \theta} \text{ Loss gradient from previous layer}$$

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

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

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

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

Loss gradient w.r.t the filter

• If you closely look at it, this represents an operation we are quite familiar with. We can represent it as a **convolution operation** between input **X** and loss gradient $\partial L / \partial \mathbf{O}$ 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

$$\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix} = \text{Input X} \quad \begin{bmatrix} \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \end{bmatrix} = \frac{\partial L}{\partial \mathbf{O}} \text{ Loss gradient from previous layer}$$

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

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

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

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

Loss gradient w.r.t the input

• If you closely look at it, this represents an operation we are quite familiar with. We can represent it as a **convolution operation between input X and loss gradient $\partial L / \partial O$ as shown below:**

For every element of X_i

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

Loss gradient w.r.t the input

• Similarly, we can expand the chain rule summation for the gradient with respect to the input. After substituting the local gradients i.e $\frac{\partial O_i}{\partial K_i}$, we have:

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

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

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

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

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

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

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

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

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

X_{11}	X_{12}	X_{13}
X_{21}	X_{22}	X_{23}
X_{31}	X_{32}	X_{33}

Input X



F_{11}	F_{12}
F_{21}	F_{22}

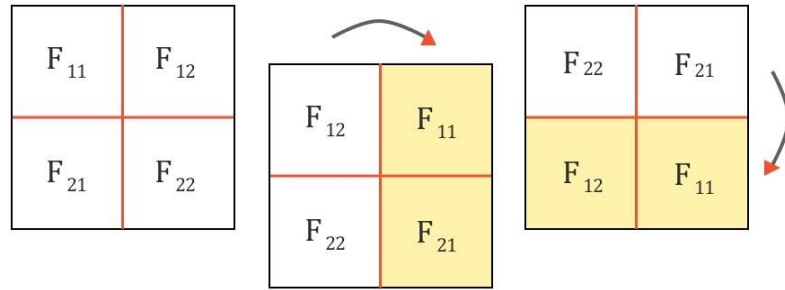
Filter F

$X_{11}F_{11}$	$X_{12}F_{12}$	X_{13}
$X_{21}F_{21}$	$X_{22}F_{22}$	X_{23}
X_{31}	X_{32}	X_{33}

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

Loss gradient w.r.t the input

- First, let us rotate the Filter F by 180 degrees. This is done by flipping it first vertically and then horizontally.



Loss gradient w.r.t the input

- We see that the loss gradient wrt the input $\frac{\partial L}{\partial X}$ is given as a full convolution between the filter and Loss gradient

F_{22}	F_{21}
F_{12}	F_{11}

Filter F

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

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

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

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

@pavisj

$\frac{\partial L}{\partial X_{11}}$	$\frac{\partial L}{\partial X_{12}}$	$\frac{\partial L}{\partial X_{13}}$
$\frac{\partial L}{\partial X_{21}}$	$\frac{\partial L}{\partial X_{22}}$	$\frac{\partial L}{\partial X_{23}}$
$\frac{\partial L}{\partial X_{31}}$	$\frac{\partial L}{\partial X_{32}}$	$\frac{\partial L}{\partial X_{33}}$

$\frac{\partial L}{\partial X}$

$$= \text{Full Convolution} \left(\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & F_{11} \\ \hline \end{array} \text{Filter F}, \begin{array}{|c|c|} \hline \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \hline \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \\ \hline \end{array} \text{Loss Gradient } \frac{\partial L}{\partial O} \right)$$

Takeaway

- Both the Forward pass and the Backpropagation of a Convolutional layer are Convolutions

$$\frac{\partial L}{\partial F} = \text{Convolution} \left(\text{Input } X, \text{ Loss gradient } \frac{\partial L}{\partial O} \right)$$

$$\frac{\partial L}{\partial X} = \text{Full Convolution} \left(\begin{array}{c} 180^\circ \text{rotated} \\ \text{Filter } F \end{array}, \text{ Loss Gradient } \frac{\partial L}{\partial O} \right)$$

Backprop in Pooling Layer

- Max Pooling
 - the error is just assigned to where it comes from
- Average Pooling
 - The error is multiplied by $1/(N \times N)$ and assigned to the whole pooling block