



# Lecture 5. Convolutional Neural Network



# Review

- 
1. Limitation of Linear Classifier
    - XOR Gate
  2. Perceptron
    - Perceptron = Linear Classifier
    - Analogy to Neurons
    - Building Block of the Neural Network
  3. MLP
    - MLP and the Neural Network
    - The Universal Approximation Theorem
    - Where Backpropagation becomes Important
  4. "Limitations of MLP"

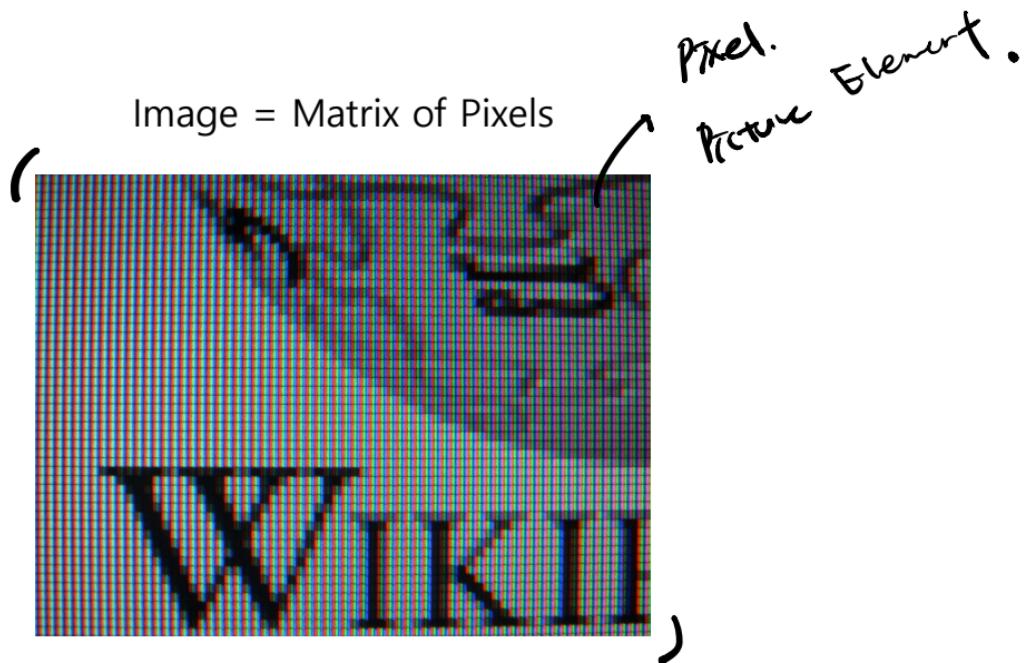
# Today's Contents

CNN

//

1. Back to Image Classification
2. Convolutional Neural Network //

# I. Back to "Image Classification" : Pixels



# Back to Image Classification : Pixels

$\partial \sim 255$

What we see



What computer sees (grayscale)

$x = [254, 081, 068, 041, 032, 071, 197,$   
 $196, 014, 132, 213, 187, 043, 041,$   
 $174, 011, 200, 254, 161, 177, 164,$   
 $201, 014, 012, 128, 242, 255, 255,$   
 $253, 212, 089, 005, 064, 196, 253,$   
 $255, 255, 251, 196, 030, 009, 165,$   
 $127, 162, 251, 254, 197, 009, 105,$   
 $062, 005, 100, 144, 097, 006, 170,$   
 $207, 083, 032, 051, 053, 134, 250]$



Image from researchgate.net

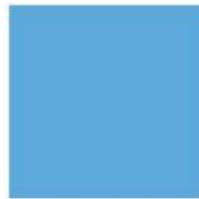
# Back to Image Classification : Pixels

$(R, G, B)$

Pixel

RGB

〃



R	G	B
61	183	228



R	G	B
255	136	73



R	G	B
105	190	40

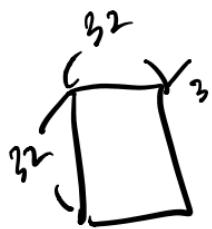
〃

Image from <https://negliadesign.com/>

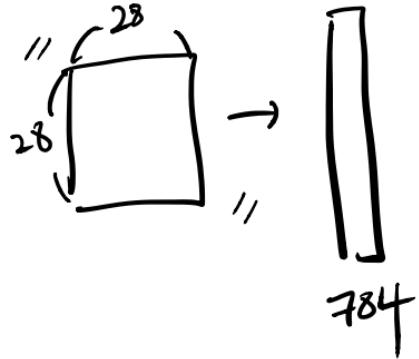
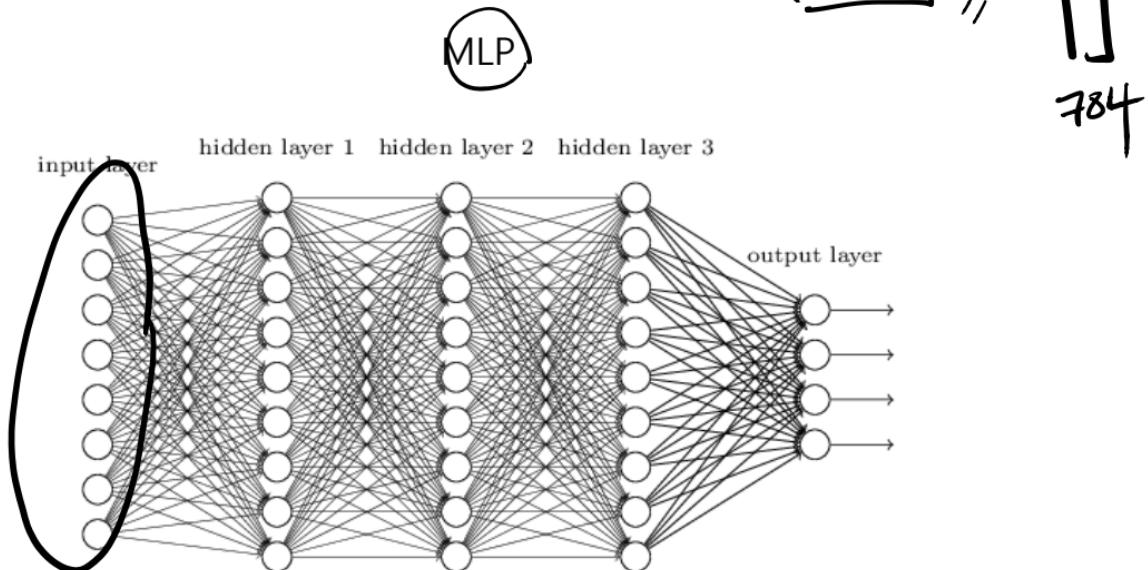


Image from <https://www.kaggle.com/>

## Back to Image Classification : MLP



$$\begin{aligned} & 32 \times 32 \times 3 \\ & = 1024 \times 3 \\ & = 3072 \end{aligned}$$



# Back to Image Classification : MLP

Must Stretch Out Tensor into Vector

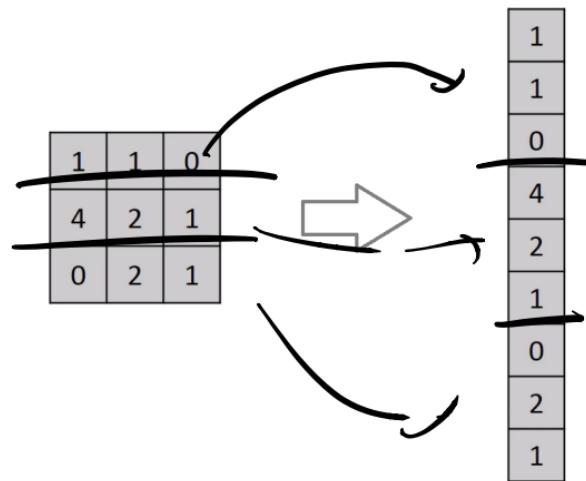


Image from <https://towardsdatascience.com>

## Back to Image Classification : MLP

Worked well on MNIST 

What about on CIFAR-10   


$$3072 \times h_1 + h_1 h_2 + h_2 \times 10$$

Back to Image Classification . MLP

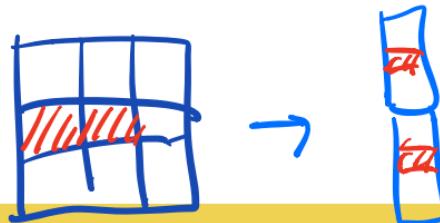
$$3072 + h_1 + h_2 + 10$$

However, MLPs are not appropriate for Image Classification...

1. Too many parameters required

2. Spatial structure is lost

↓ Matrix Tensor. → Vector.



## Back to Image Classification : MLP

We want a model that,

1. Uses fewer parameters than MLPs
2. Preserves the spatial structure

## Back to Image Classification : MLP



“Recall The Universal Approximation Theorem...”

“( “1개의 Hidden Layer를 가진 MLP로 어떤 함수도 근사할 수 있다.” )

However,

“( “A feedforward network with a single layer is sufficient to represent any function,  
But the layer may be infeasibly large and may fail to learn and generalize correctly.” )

# CNN : Idea

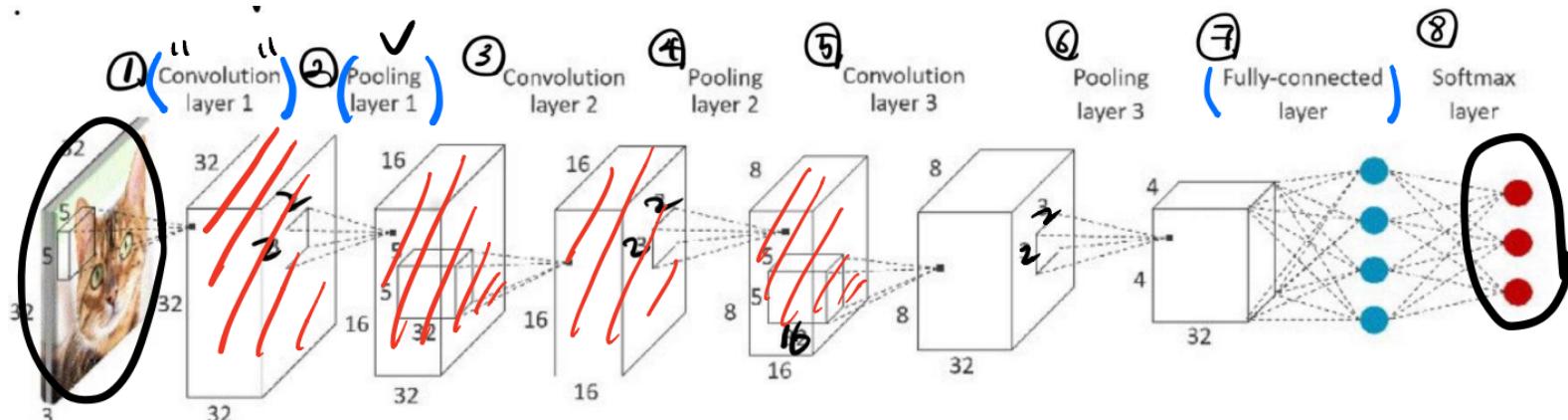
We want a model that,

- // 1. Uses fewer parameters than MLPs
- 2. Preserves the spatial structure //

**Convolutional Neural Network** to the rescue!!!

# CNN : Idea

Image from <https://community.arm.com/>



What happened?!?!?!

# CNN : Idea

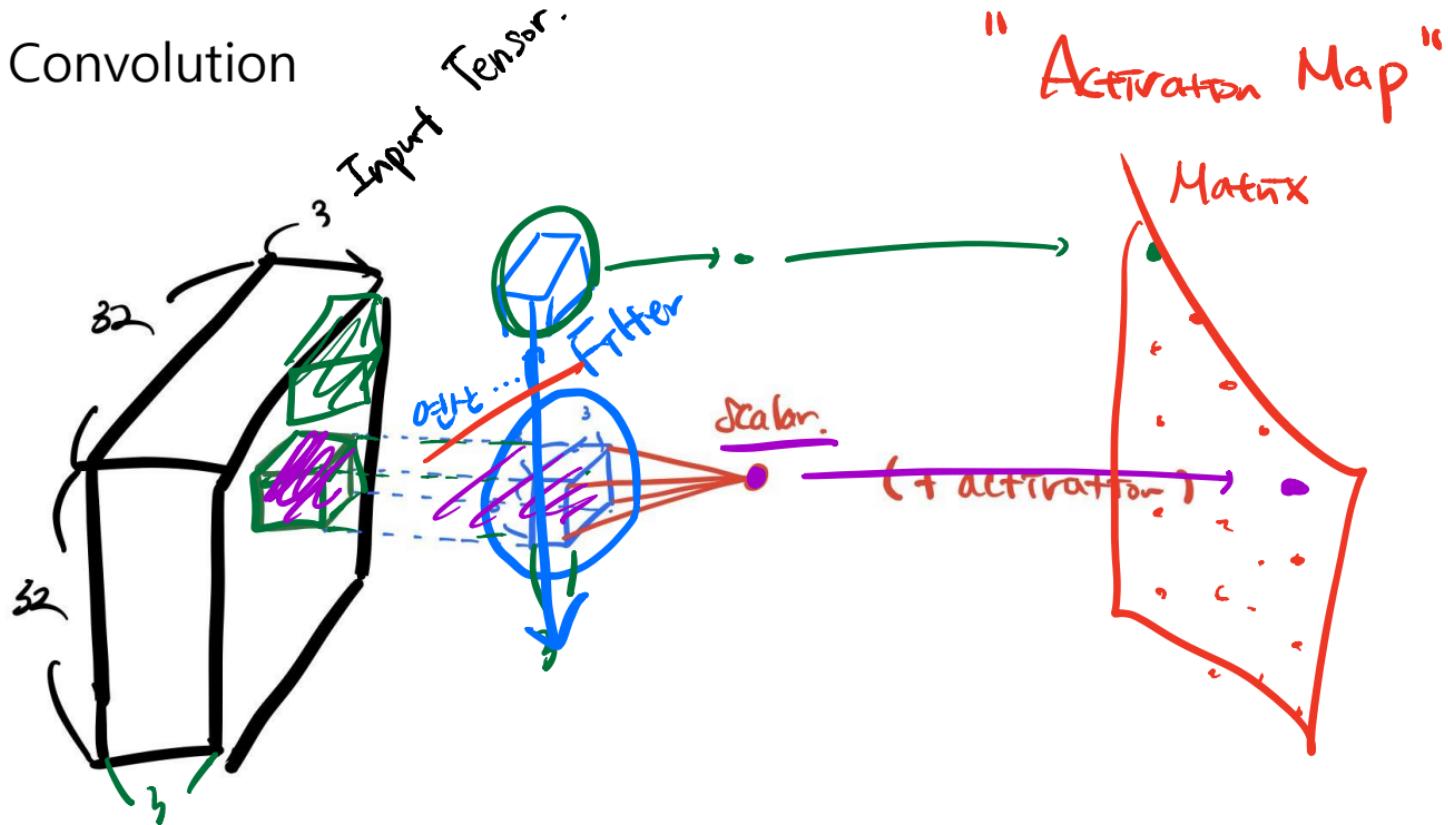
- 1 What happens during Tensor → Tensor Mapping?
- 2 What are Convolution, Pooling, Fully-Connected Layers?
- 3 How does this Model extract features from spatial structure?
- 4 Does it use less parameters than the MLP?

# CNN : Convolution

How do we map Tensor to Tensor, while preserving the spatial structure?

Apply **Filter** to small regions in the **Tensor** (Image)

## CNN : Convolution

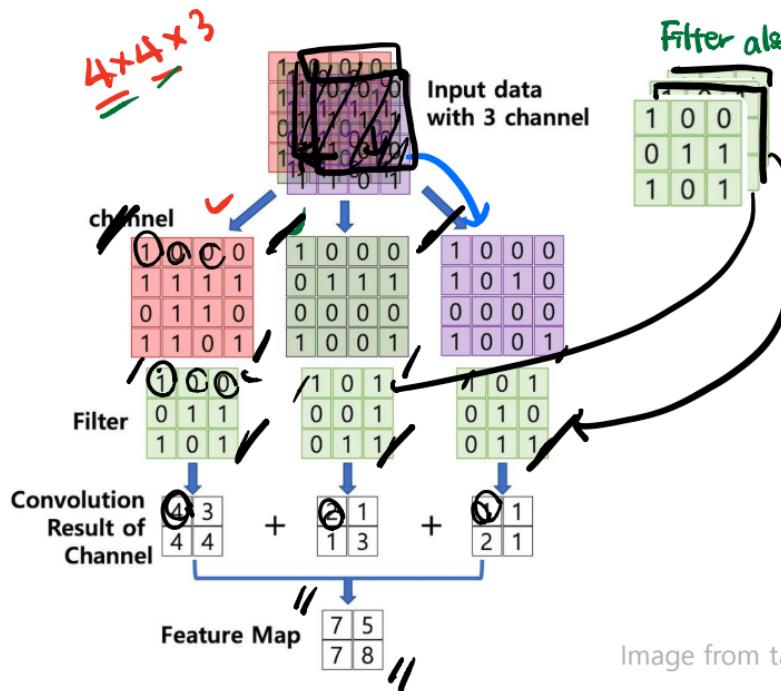


$$( + 0 + 0 + 0 + 1 + 1 + 0 + 0 + 1 ) = 4$$

## CNN : Convolution

$3 \times 3 \times 3$

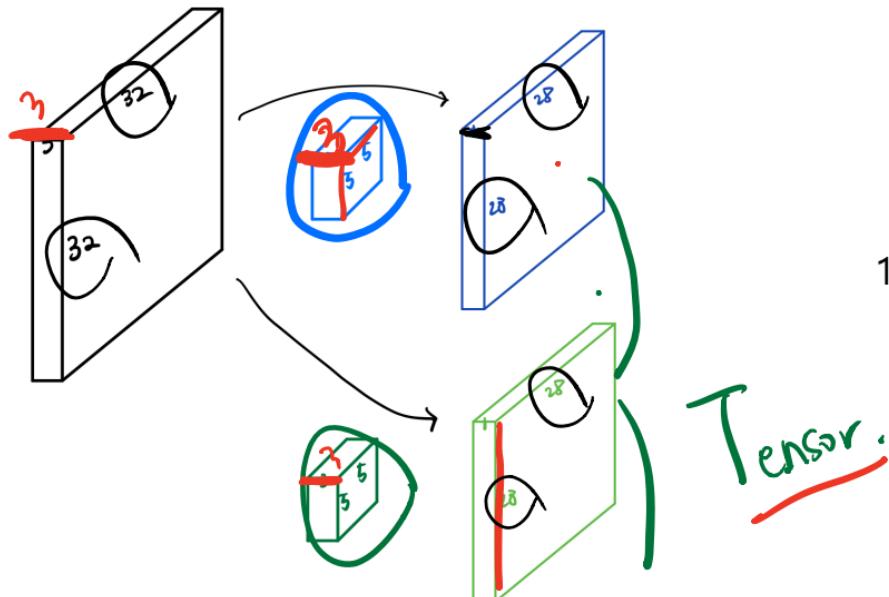
ReLU  $y = \max(0, x)$



While sliding over the tensor,  
calculate the  
sum of elementwise multiplication  
This is called, "Convolution"

Image from taewan.kim/ppst/cnn

# CNN : Convolution



$$T \rightarrow M$$

1 Filter = 1 Activation Map

T

# CNN : Convolution

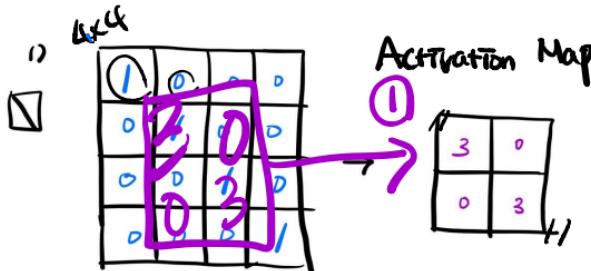
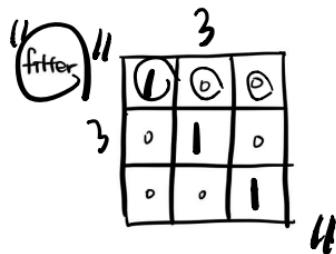
MLP X .

Then, how do filters extract features from the image,  
while preserving the spatial structure?

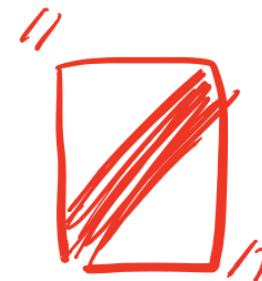
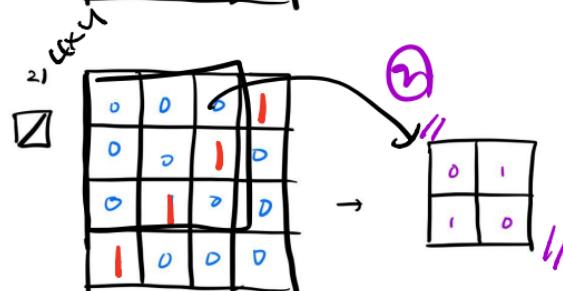
# CNN : Convolution



For simplicity, assume White = 0 , Black = 1, and activation = RELU



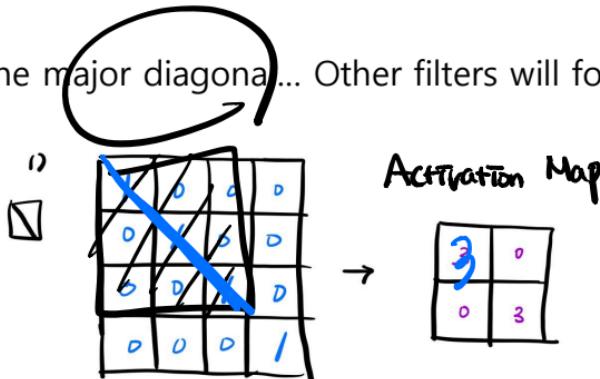
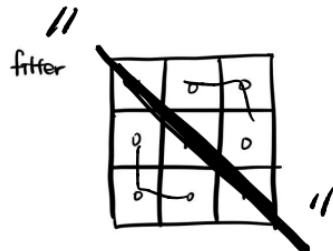
Activation Map



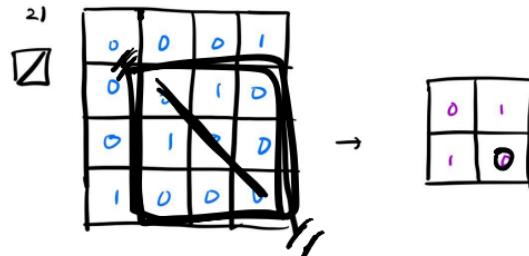
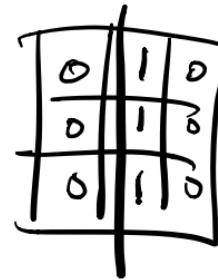
# CNN : Convolution



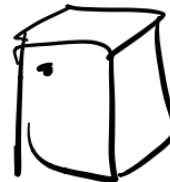
This filter focuses on the major diagonal... Other filters will focus on other features!



Activation Map



# CNN : Convolution



Filters extract features from the spatial structure of the tensor

" " "  
They are Parameters!!!

" " "  
\*\* So, they must be Updated via Backpropagation.

For more detail, read <https://becominghuman.ai/back-propagation-in-convolutional-neural-networks-intuition-and-code-714ef1c38199>

I O F S

CNN : Convolution

$$O = \frac{I - F}{S} + 1$$

$$S=1.$$
$$I=5. F=3$$

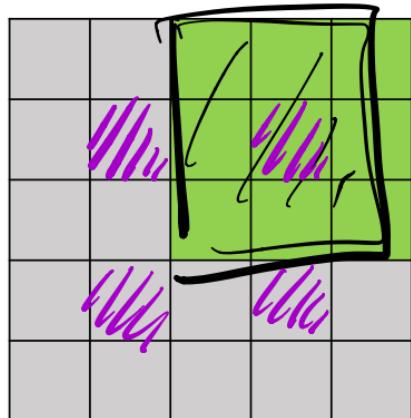
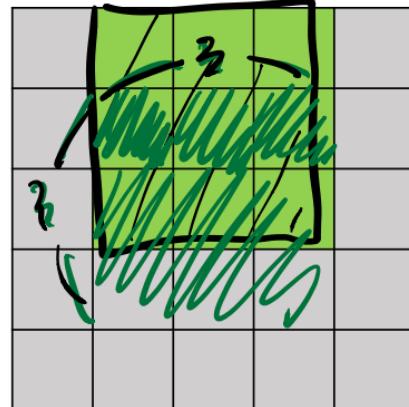
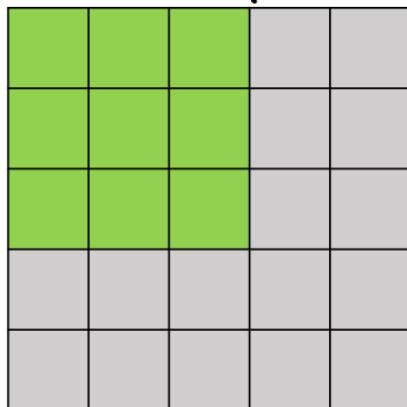
How much does a filter move per slide? **Stride**

$$\frac{5-3}{1} + 1 = 3$$



Stride=1

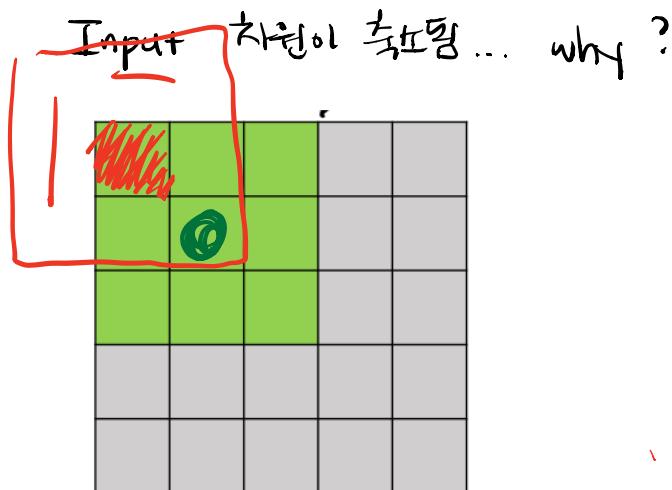
Stride=2



CNN: Padding.

~~Fix~~

$$O = \frac{I - F}{S} + 1 \leq I - F + 1 < I$$



Remedy this via, "Padding"

The diagram shows a 3x3 input grid with a 2x2 kernel. The output grid has a green circle at its center. A red box highlights the top-left 2x2 subgrid of the input, which is processed by the kernel to produce the output value. The input grid has zeros added around the edges, demonstrating padding.

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0



However,

Data Log  $\leftrightarrow$  Notse.

## CNN : Convolution



$$O_w = \frac{I_w + 2P_{aw} - F_w}{S_w} + 1$$

O = output size

N = input size

F = filter size

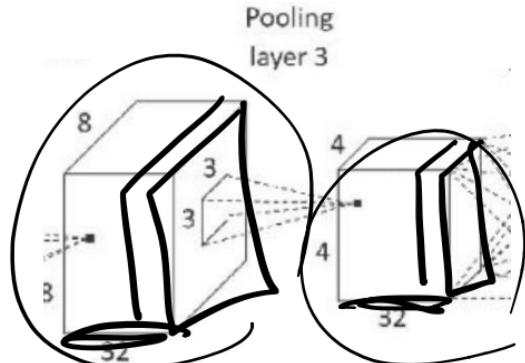
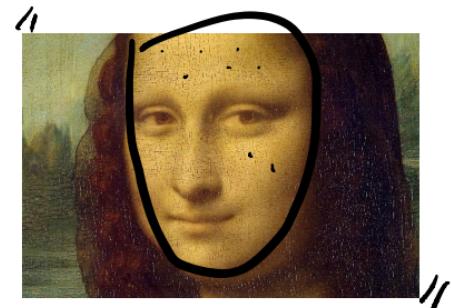
S = stride

Pa = padding size

# CNN : Pooling

We want to reduce the dimension of the feature,

via **Pooling**



# CNN : Pooling

Extract the representative value over a region... "No parameter involved!!!"

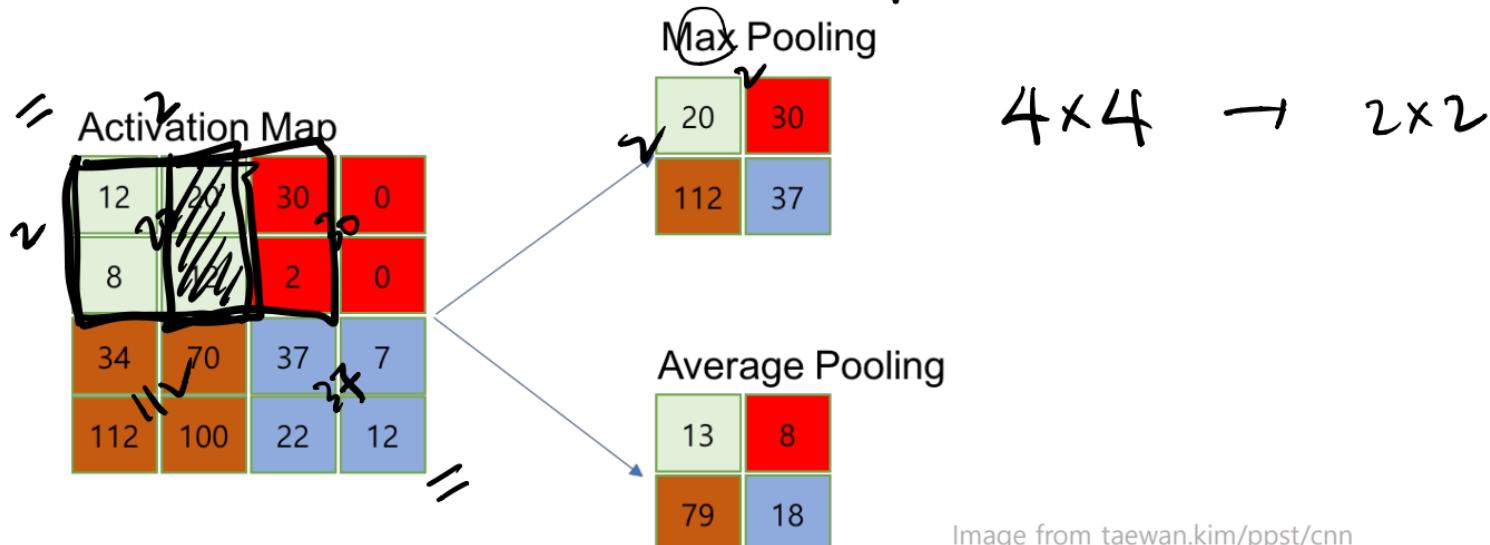


Image from taewan.kim/ppst/cnn

## CNN : Pooling

$$O = \frac{I}{P_o}$$

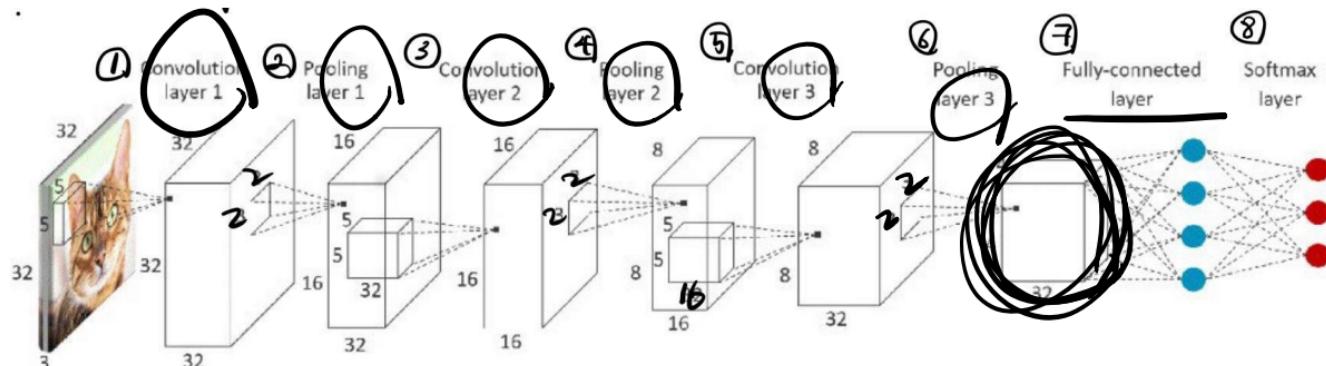
O = output size

I = input size

Po = pool size

\*\* Generally, Pool Size = Pooling Stride

# CNN : Fully-Connected Layer



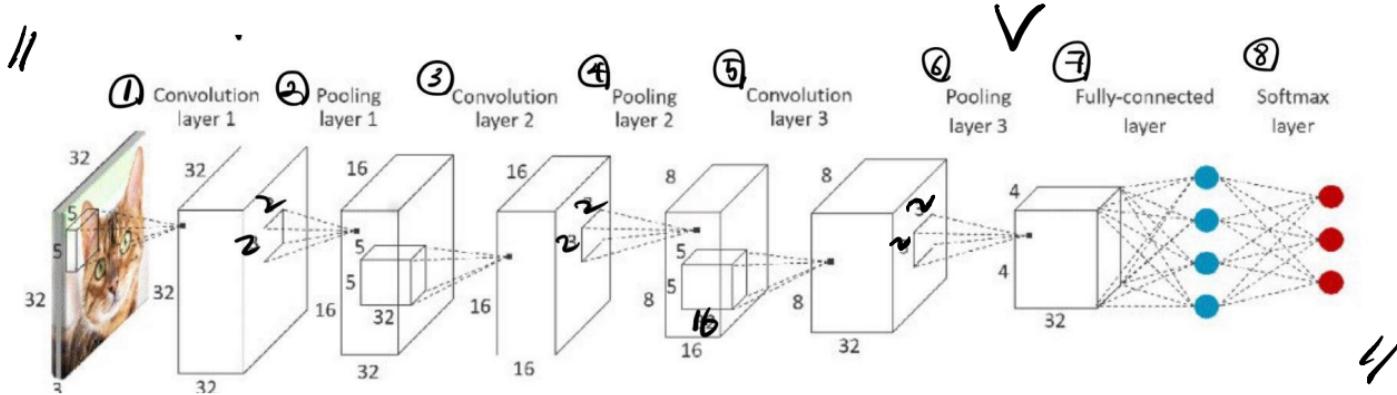
Extract features with **CONV Layers**

Then, reduce dimension with **POOL Layers**

Finally, flatten data and make classification via **FC(Fully-Connected) Layer**

# CNN : Architecture

Conv + Pool ✓



Extract features with **CONV Layers**

Then, reduce dimension with **POOL Layers**

Finally, flatten data and make classification via **FC(Fully-Connected) Layer**

# CNN : Architecture

Does it use less parameters than the MLP?

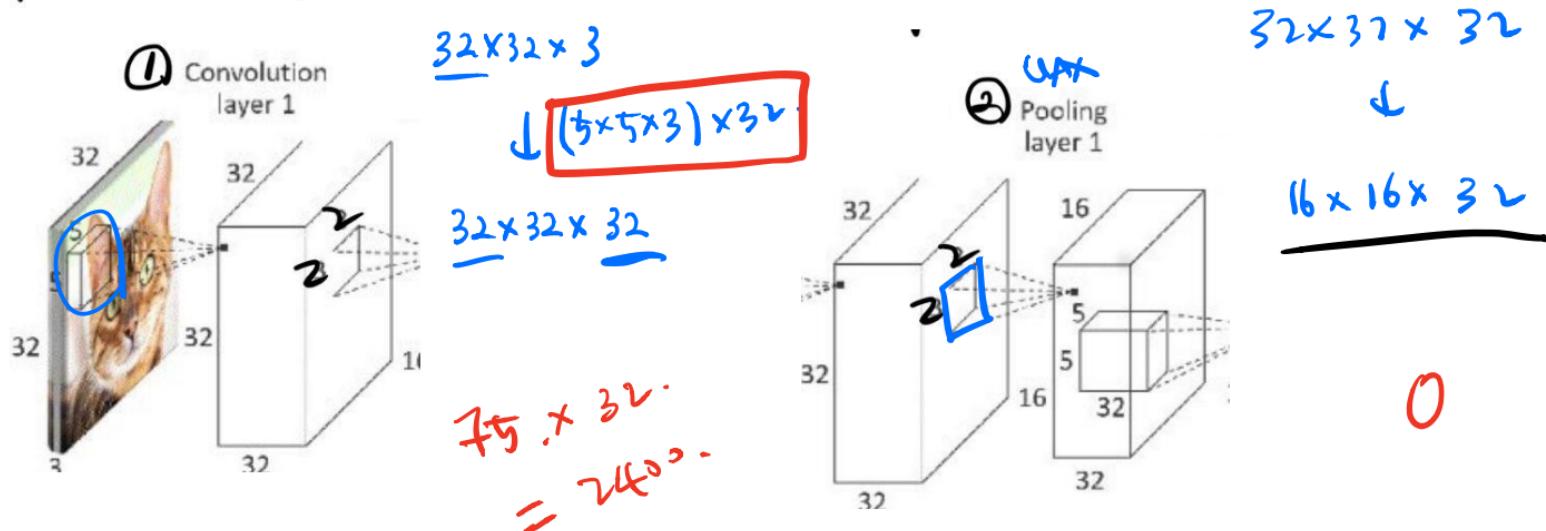
Let's follow the CNN Architecture and figure it out!!!

# CNN : Architecture

CONV1

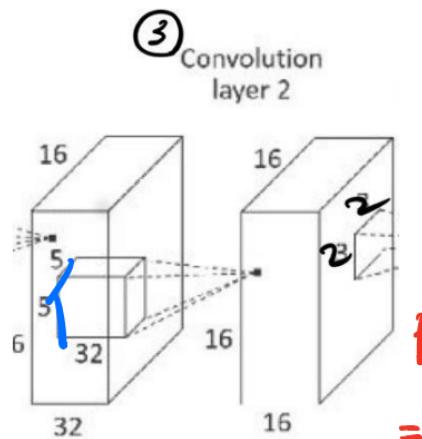
$$\frac{32 + 2 \times 2 - 5}{1} + 1 = \underline{\underline{32}}$$

POOL1

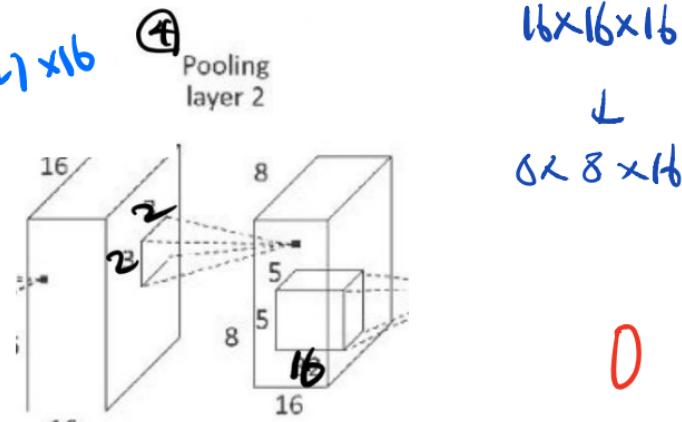


# CNN : Architecture

CONV2

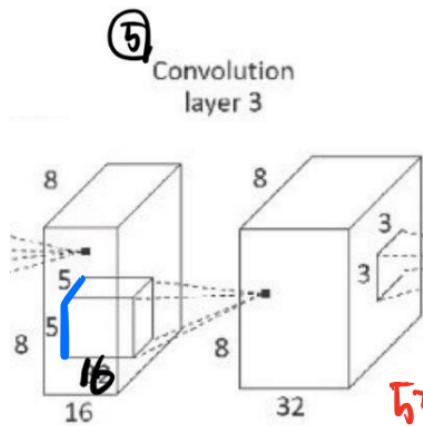


MAX  
POOL2

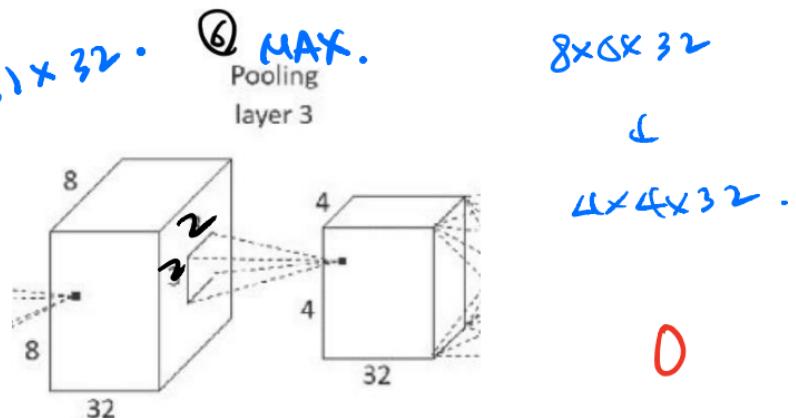


# CNN : Architecture

CONV3

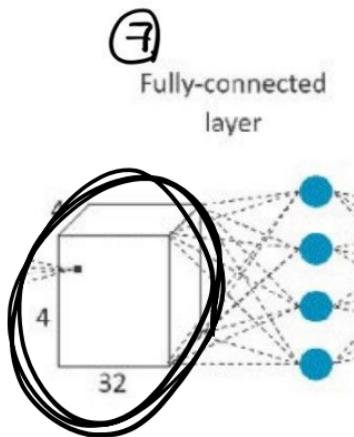


POOL3



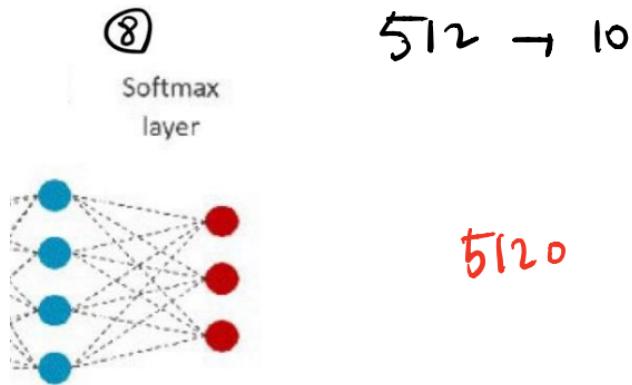
# CNN : Architecture

FC



$4 \times 4 \times 32$   
↓  
0  
512

Softmax



## CNN : Architecture

$$2400 + 0 + 12800 + 0 + 12800 + 0 + 0 + 5120 = 33120.$$

Total ( 33120 . ) Parameters were Used!

## CNN : Architecture

$$\begin{aligned} & 3072^+ \xrightarrow{(16 \times 16 \times 32)} \\ & + 16 \times 16 \times 32 + 8 \times 8 \times 16 \\ & + \dots \\ & = H. \end{aligned}$$

However, if we used MLP instead, assuming CONV + POOL as one Hidden Layer,

$$\begin{aligned} & 3072 \\ & \cancel{32 * 32 * 3} \rightarrow \cancel{16 * 16 * 32} \rightarrow \cancel{8 * 8 * 16} \rightarrow \cancel{4 * 4 * 32} \rightarrow 10 \end{aligned}$$

Total ( 35107840 ) Parameters were Used!

# CNN : Architecture

CNN uses much less Parameters than the MLP!!!

# CNN : Performance

CIFAR-10			
Who is the best in CIFAR-10 ?			
Result	Method	Venue	Details
96.53%	Fractional Max-Pooling	arXiv 2015	<a href="#">Details</a>
95.59%	Striving for Simplicity: The All Convolutional Net	ICLR 2015	<a href="#">Details</a>
94.16%	All you need is a good init	ICLR 2016	<a href="#">Details</a>
94%	Lessons learned from manually classifying CIFAR-10	unpublished 2011	<a href="#">Details</a>
93.95%	Generalizing Pooling Functions in Convolutional Neural Networks: Mixed, Gated, and Tree	AISTATS 2016	<a href="#">Details</a>
93.72%	Spatially-sparse convolutional neural networks	arXiv 2014	<a href="#">Details</a>
93.63%	Scalable Bayesian Optimization Using Deep Neural Networks	ICML 2015	<a href="#">Details</a>
93.57%	Deep Residual Learning for Image Recognition	arXiv 2015	<a href="#">Details</a>
93.46%	Fast and Accurate Deep Network Learning by Exponential Linear Units	arXiv 2015	<a href="#">Details</a>
93.34%	Universum Prescription: Regularization using Unlabeled Data	arXiv 2015	<a href="#">Details</a>
93.25%	Batch-normalized Maxout Network in Network	arXiv 2015	<a href="#">Details</a>
93.13%	Competitive Multi-scale Convolution	arXiv 2015	<a href="#">Details</a>
92.91%	Recurrent Convolutional Neural Network for Object Recognition	CVPR 2015	<a href="#">Details</a>
92.49%	Learning Activation Functions to Improve Deep Neural Networks	ICLR 2015	<a href="#">Details</a>

96.53%

30%

Classification Rankings :  
[https://rodrigob.github.io/are\\_we\\_there\\_yet/build/classification\\_datasets\\_results.html](https://rodrigob.github.io/are_we_there_yet/build/classification_datasets_results.html)

CNN : Performance

50%.



# Review

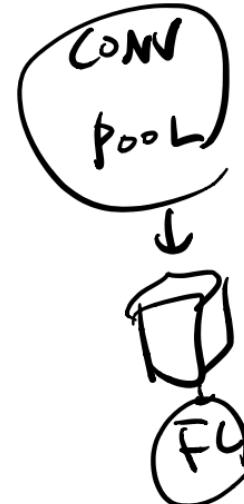
MLP

We want a model that,

1. Uses fewer parameters than MLPs
2. Preserves the spatial information

CNN

**Convolutional Neural Network** to the rescue!!!



# Preview on Next Lecture(s)



## Limitations of MLP

1. Needs a lot of Labeled Data
2. Vanishing Gradient
3. Overfitting
4. Gets stuck in the Local Minima / Saddle Point



# Preview on Next Lecture(s)

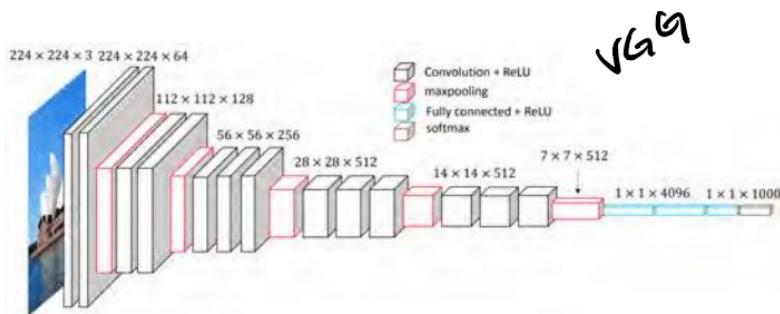


Image from researchgate.net

VGG19



Image from medium.com

GoogLeNet