

SK hynix i-TAP 반도체 Data Scientist를 위한 ML/DL 심화 커리큘럼

1강. Convolutional Neural Networks

Ernest K. Ryu (류경석)

2020.11.06

1강. Outline

- Basics of (non-adaptive) SGD
- PyTorch as a GPU-computing numerical library
- Backpropagation
- Multilayer perceptron
- Convolutional neural networks

Optimization

We consider

$$\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

where f_1, \dots, f_N are “differentiable” functions.

In DL, the ReLU activation function $\sigma(x) = \max(0, x)$ is said to be “differentiable”.

Gradient Descent (GD)

Define $F(\theta) = \frac{1}{N} \sum_{i=1}^N f_i(\theta)$. Then
$$\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \quad F(\theta)$$

GD:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla F(\theta^k)$$

where $\alpha_0, \alpha_1, \dots \in \mathbb{R}$ are stepsizes.

Since $\nabla F = \frac{1}{N} \sum_{i=1}^N \nabla f_i$, can parallelize $\nabla F(\theta^k)$ computation.

Why does GD converge?

Taylor expansion of F about θ^k :

$$F(\theta) = F(\theta^k) + \nabla F(\theta^k)^T (\theta - \theta^k) + \mathcal{O}\left((\theta - \theta^k)^2\right)$$

Plug in θ^{k+1} :

$$F(\theta^{k+1}) = F(\theta^k) - \alpha_k \|\nabla F(\theta^k)\|^2 + \mathcal{O}(\alpha_k^2)$$

$-\nabla F(\theta^k)$ is steepest descent direction. For small (cautious) α_k , GD step reduces function value.

Stochastic Gradient Descent (SGD)

We consider

$$\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{N} \sum_{i=1}^N f_i(\theta) = \mathbb{E}_{I \sim \text{Uniform}\{1, \dots, N\}} [f_I(\theta)]$$

SGD:

$$\begin{aligned} i(k) &\sim \text{Uniform}\{1, \dots, N\} \\ \theta^{k+1} &= \theta^k - \alpha_k \nabla f_{i(k)}(\theta^k) \end{aligned}$$

where $i(k)$ are uniform IID indices.

Why does SGD converge?

$\nabla f_{i(k)}(\theta^k)$ is an unbiased estimate of the gradient $\nabla F(\theta^k)$

$$\mathbb{E}_{i(k)} \nabla f_{i(k)}(\theta^k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\theta^k) = \nabla F(\theta^k)$$

(So $\nabla f_{i(k)}(\theta^k)$ is a “stochastic gradient” of F at θ^k)

Why does SGD converge?

Plug θ^{k+1} into Taylor expansion of F about θ^k :

$$F(\theta^{k+1}) = F(\theta^k) - \alpha_k \nabla F(\theta^k)^T \nabla f_{i(k)}(\theta^k) + \mathcal{O}(\alpha_k^2)$$

Expectation on both sides:

$$\mathbb{E}_k F(\theta^{k+1}) = F(\theta^k) - \alpha_k \|\nabla F(\theta^k)\|^2 + \mathcal{O}(\alpha_k^2)$$

(\mathbb{E}_k is expectation conditioned on θ^k)

$-\nabla f_{i(k)}(\theta^k)$ is descent direction in expectation. For small (cautious) α_k , SGD step reduces function value in expectation.

SGD with general expectation

Consider general expectation

$$\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \quad \mathbb{E}_{\omega} [f_{\omega}(\theta)]$$

where ω is a random variable. Expectation (not finite sum) appears when you have generative model of ω . E.g. GAN.

SGD:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{\omega^k}(\theta^k)$$

where $\omega^0, \omega^1, \dots$ IID random samples of ω .

Mini-batch SGD

For each k , let $i(k, 1), \dots, i(k, B)$ be IID indices.

$$\frac{1}{B} \sum_{j=1}^B \nabla f_{i(k,j)}(\theta^k)$$

is also an unbiased estimate of $\nabla F(\theta^k)$ since

$$\mathbb{E} \frac{1}{B} \sum_{j=1}^B \nabla f_{i(k,j)}(\theta^k) = \frac{1}{B} \sum_{j=1}^B \mathbb{E} \nabla f_{i(k,j)}(\theta^k) = \frac{1}{B} \sum_{j=1}^B \nabla F(\theta^k) = \nabla F(\theta^k)$$

Mini-batch SGD

Mini-batch SGD:

$$g = 0$$

For $j = 1, \dots, B$

$$i(k, j) \sim \text{Uniform}\{1, \dots, N\}$$

$$g = g + \frac{1}{B} \nabla f_{i(k, j)}$$

$$\theta^{k+1} = \theta^k - \alpha_k g$$

is also an instance of SGD.

Mini-batch Size

Mathematically (measuring performance per iteration)

- Use large batch is when noise/randomness is large.
- Use small batch is when noise/randomness is small.

Practically (measuring performance per unit time)

- Large batch allows more efficient communication and computation, up to the GPU memory limit.
- Often best to increase batch size up to the GPU memory limit.

Cyclic (mini-batch) SGD

Cyclic SGD:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{i(k)}(\theta^k)$$

where $i(k)$ is selected in a cyclic order.

Can also write as:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{\text{mod}(k,N)+1}(\theta^k)$$

Strictly speaking, is not an instance of SGD as unbiased estimation property lost.

Epoch in Optimization and Training

Epoch: loosely defined as computation time of computing 1 full gradient. One iteration of GD is, by definition, an epoch. N iterations of SGD constitute an epoch.

Epoch is a convenient unit for counting iterations (rather than directly counting iteration numbers).

Cyclic (mini-batch) SGD

Cyclic SGD advantage:

- Simpler than SGD
- Uses all datapoints within single epoch.

Cyclic SGD disadvantage:

- Worse than SGD in some cases, theoretically and empirically.
- In deep learning, neural nets learn to anticipate cyclic order.

Shuffled Cyclic (mini-batch) SGD

Shuffled cyclic SGD:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{i(k)}(\theta^k)$$

where $i(k)$ is selected in a cyclic order shuffled every epoch.

Can also write as:

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{\sigma\left\lfloor \frac{k}{N} \right\rfloor (\bmod(k, N) + 1)}(\theta^k)$$

where $\sigma^0, \sigma^1, \dots$ is a sequence of random permutations.

Shuffled Cyclic (mini-batch) SGD

Shuffled Cyclic SGD:

For $e = 1, \dots, E$ *//for each epoch*

$\sigma \sim \text{randomPermutation}(N)$

For $i = 1, \dots, N$

$$\theta^{k+1} = \theta^k - \alpha_k \nabla f_{\sigma(k)}(\theta^k)$$

$$k = k + 1$$

Shuffled Cyclic (mini-batch) SGD

Shuffled cyclic SGD advantage:

- Uses all datapoints within single epoch.
- Network cannot learn to anticipate data order.
- Generally best performance.

Shuffled cyclic SGD advantage:

- Not as simple. (But PyTorch makes it simple to use.)
- Theory not as strong as regular SGD.

PyTorch: GPU Numerical Computing Library

PyTorch is a machine learning library of Python, but PyTorch is fundamentally a numerical computation library.

Features of PyTorch that make it suitable for using neural networks and machine learning:

- PyTorch supports easy GPU computation.
- Automatic differentiation.
- Numerous ML libraries and sample code.

GPU Computing Code With CUDA

GPU computing operations:

- cudaMemcpy (CPU→GPU or GPU→CPU)
- CPU code
- GPU kernel calls (CPU instructs GPU to execute computation)

GPU computing workflow:

- (i) send data CPU→GPU
- (ii) Compute on GPU
- (iii) Receive result GPU→CPU

CPU vs. GPU variables

Variables either reside in CPU or GPU memory.

- CPU variable computation on CPU
- GPU variable computation on GPU

CPU and GPU variables cannot directly interact. Can interact only after CPU→GPU or GPU→CPU transfer.

PyTorch Demo

Power iteration example on PyTorch

```
send A from host (CPU) to device (GPU)
send x=x0 from host (CPU) to device (GPU)
for _ in range(100):
    tell GPU to compute  $x=A*x$ 
send x from device (GPU) to host (CPU)
```

Back Propagation \subseteq Automatic Differentiation

Automates gradient computation! Only need to specify how to evaluate function.

Gradient costs roughly $5 \times$ computation cost* of evaluating function.

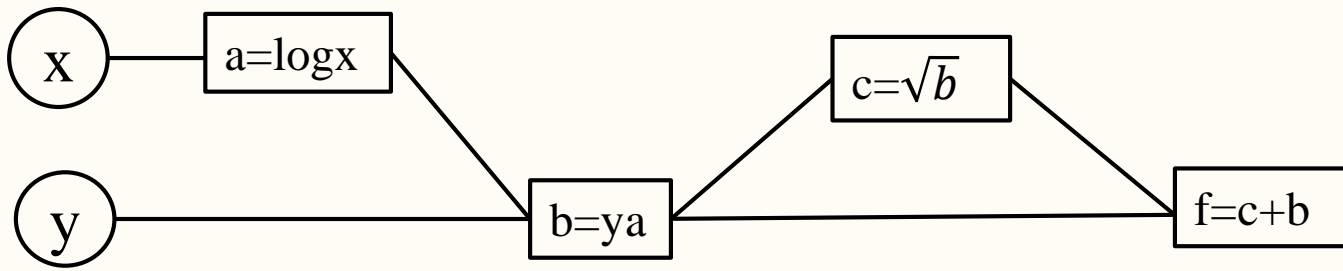
AutoDiff is not

- Finite differencing
- Symbolic differentiation.

AutoDiff \approx chain rule of vector calculus

Chain Rule

- Consider $f(x, y) = y \log x + \sqrt{y \log x}$
- Evaluate f with the computation graph:



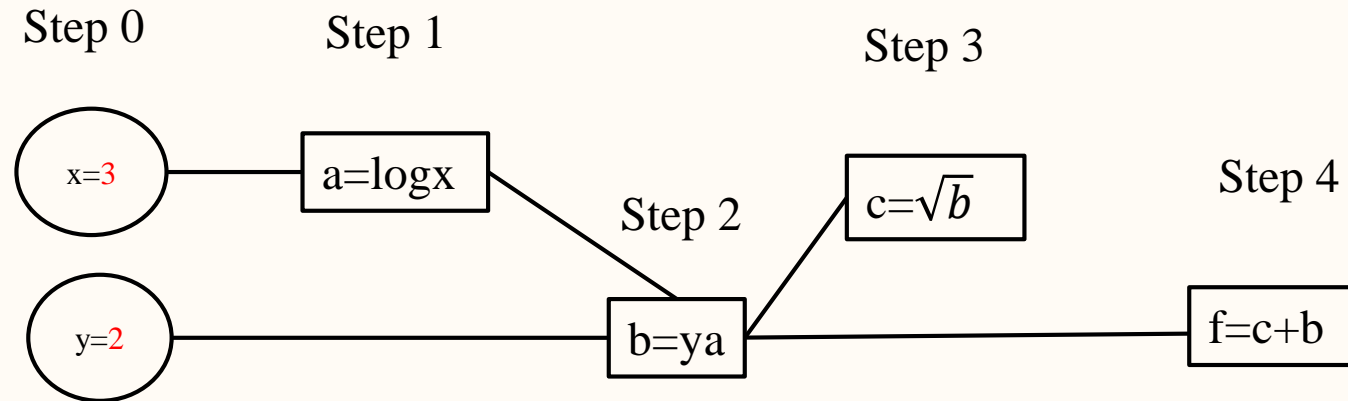
- Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} \left(\frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \frac{\partial x}{\partial x} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial x} \right) + \frac{\partial f}{\partial b} \left(\frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \frac{\partial x}{\partial x} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial x} \right)$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} \left(\frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \frac{\partial x}{\partial y} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial y} \right) + \frac{\partial f}{\partial b} \left(\frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \frac{\partial x}{\partial y} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial y} \right)$$

But in what order do you evaluate the chain rule expression?

Forward mode auto-diff



$$0. x = 3, y = 2, \frac{\partial x}{\partial x} = 1, \frac{\partial x}{\partial y} = 0, \frac{\partial y}{\partial x} = 0, \frac{\partial y}{\partial y} = 1$$

$$1. a = \log x = \log 3, \frac{\partial a}{\partial x} = \frac{1}{x} \cdot \frac{\partial x}{\partial x}, \frac{\partial a}{\partial y} = 0$$

$$2. b = ya = 2\log 3, \frac{\partial b}{\partial x} = \frac{\partial y}{\partial x} a + y \frac{\partial a}{\partial x} = \frac{2}{3}, \frac{\partial b}{\partial y} = \frac{\partial y}{\partial y} a + y \frac{\partial a}{\partial y} = a = \log 3$$

Computation does not involve 'x' or derivatives of 'x'

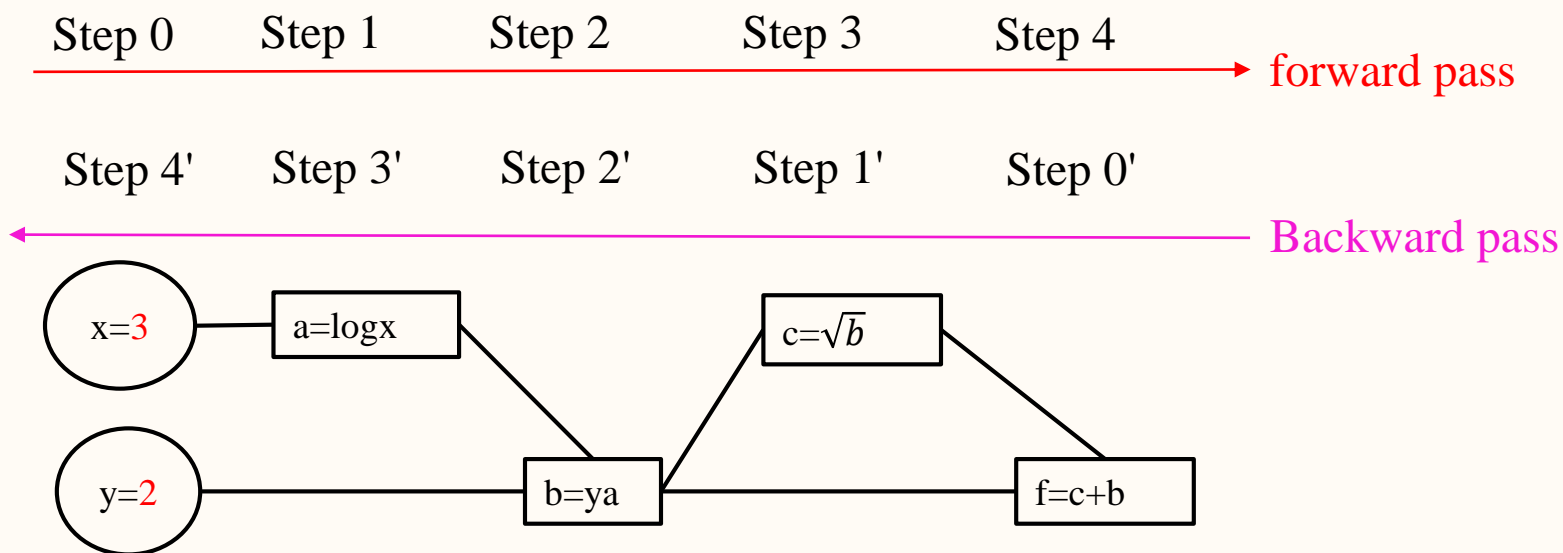
$$3. c = \sqrt{b} = \sqrt{2\log 3}, \frac{\partial c}{\partial x} = \frac{1}{\sqrt{b}} \frac{\partial b}{\partial x} = \frac{2}{3\sqrt{2\log 3}}, \frac{\partial c}{\partial y} = \frac{1}{\sqrt{b}} \frac{\partial b}{\partial y} = \sqrt{\frac{\log 3}{2}}$$

Computation only depends on node 'b'

$$4. f = c + b = \sqrt{2\log 3} + 2\log 3, \frac{\partial f}{\partial x} = \frac{\partial c}{\partial x} + \frac{\partial b}{\partial x} = \frac{2}{3} \left(1 + \frac{1}{\sqrt{2\log 3}} \right), \frac{\partial f}{\partial y} = \frac{\partial c}{\partial y} + \frac{\partial b}{\partial y} = \sqrt{\frac{\log 3}{2}} + \log 3$$

Computation only depends on node 'b' and 'c'

Reverse mode auto-diff (Backpropagation)



0. $x = 3, y = 2$
1. $a = \log 3$
2. $b = 2 \log 3$
3. $c = \sqrt{2 \log 3}$
4. $f = \sqrt{2 \log 3} + 2 \log 3$

$$\begin{aligned}
 0'. \quad \frac{\partial f}{\partial f} &= 1 \\
 1'. \quad \frac{\partial f}{\partial c} &= \frac{\partial f}{\partial f} \frac{\partial f}{\partial c} = 1 \\
 2'. \quad \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial f}{\partial f} \frac{\partial f}{\partial b} = \frac{1}{\sqrt{b}} + 1 = \frac{1}{\sqrt{2 \log 3}} + 1 \\
 3'. \quad \frac{\partial f}{\partial a} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial a} = 2 \left(\frac{1}{\sqrt{2 \log 3}} + 1 \right) \\
 4'. \quad \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} = \frac{2}{3} \left(\frac{1}{\sqrt{2 \log 3}} + 1 \right) \\
 \frac{\partial f}{\partial y} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial y} = \left(\frac{1}{\sqrt{2 \log 3}} + 1 \right) a = \sqrt{\frac{\log 3}{2}} + \log 3
 \end{aligned}$$

Backward pass depends on node values computed in forward pass.

Autodiff by Jacobian multiplication

Consider $g = f_1 \circ f_2 \circ \cdots \circ f_N$ where $f_i: R^{n_i} \rightarrow R^{n_{i-1}}$ for $i = 1, \dots, N$.

Chain rule: $\nabla_x g(x) = Df_1 \quad Df_2 \quad \cdots \quad Df_N$ **D denotes Jacobian.**
 $n_0 \times n_1 \quad n_1 \times n_2 \quad n_{N-1} \times n_N$

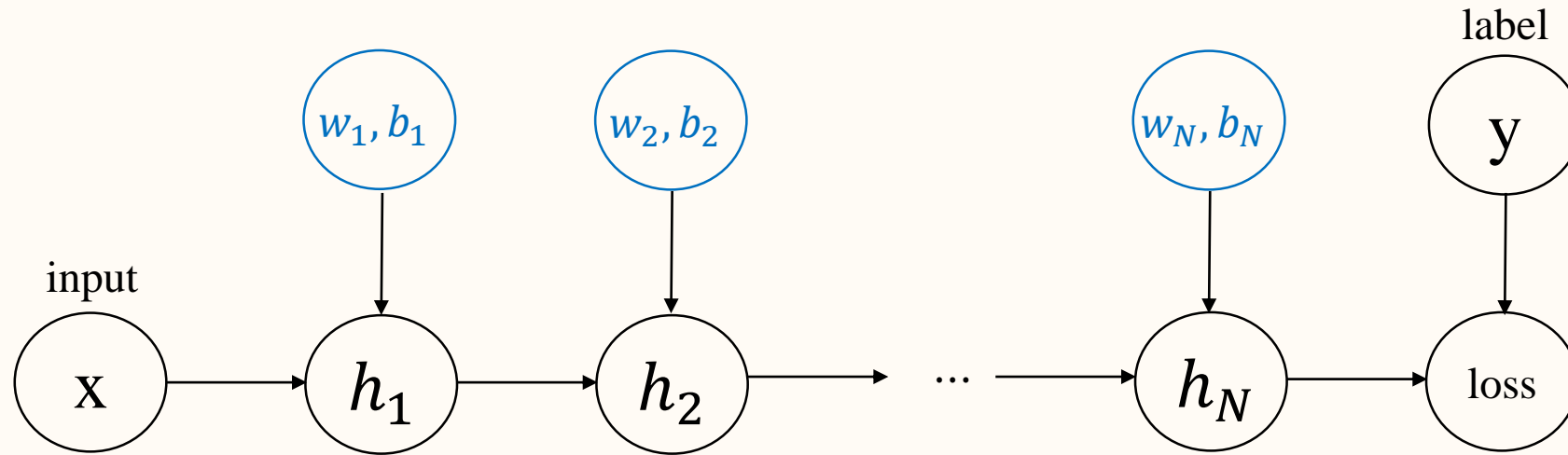
Forward mode: $Df_1(Df_2(\cdots(Df_{N-1}Df_N) \cdots))$

Reverse mode: $((Df_1 Df_2) Df_3) \cdots) Df_N$



Optimal if $n_0 \leq n_1 \leq \cdots \leq n_N$.
Proof by dynamic programming.

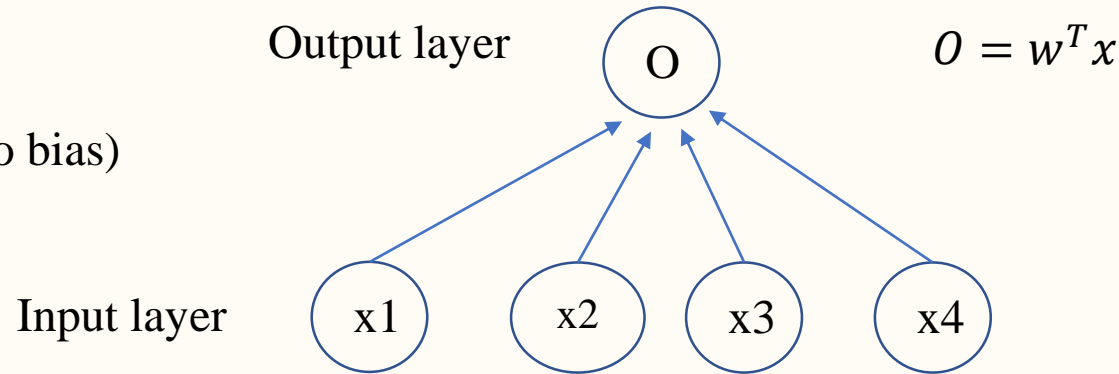
Backprop on multi-layer perceptron



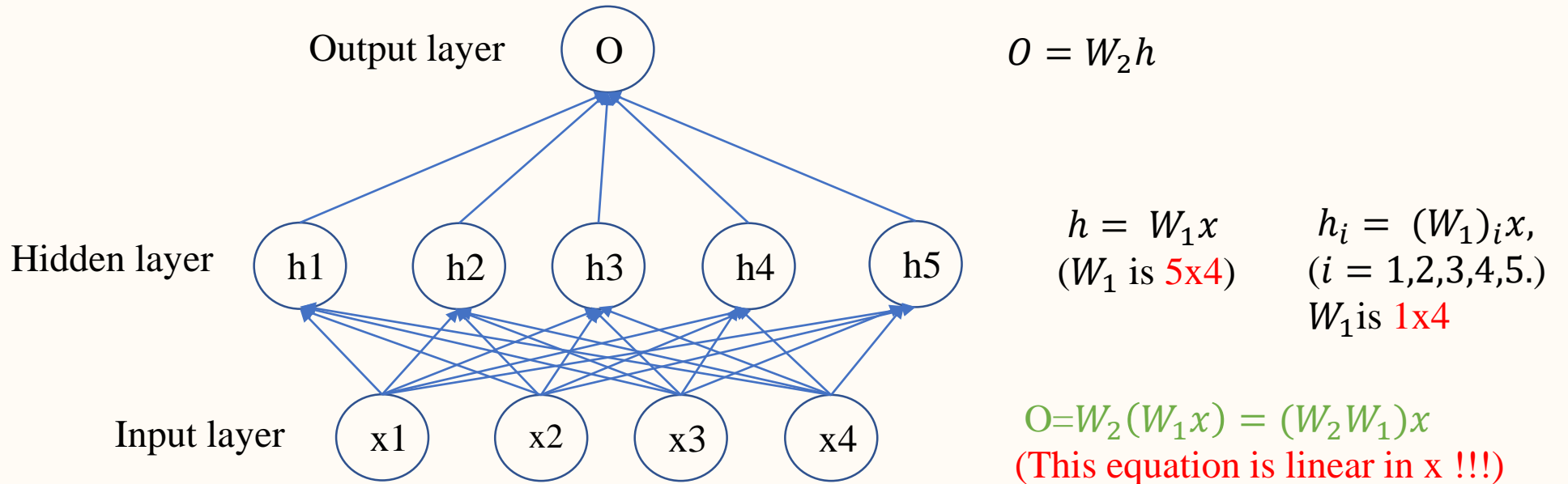
- In NN training, [parameters](#) and fixed inputs are distinguished.
 - In Pytorch, 1. evaluate the loss function
2. call `backward()` to perform backward pass and compute gradients.
-
- When performing the forward pass, intermediate node values are stored so that they can later be used in backward pass.
 - In testing loop, we don't compute gradients so this is unnecessary.
 - The [torch.no_grad\(\) context manager](#) allows intermediate node values to be discarded or not be stored. This saves memory and can reduce the time to compute the test loop.

Multilayer Perceptron (fully connected deep neural network)

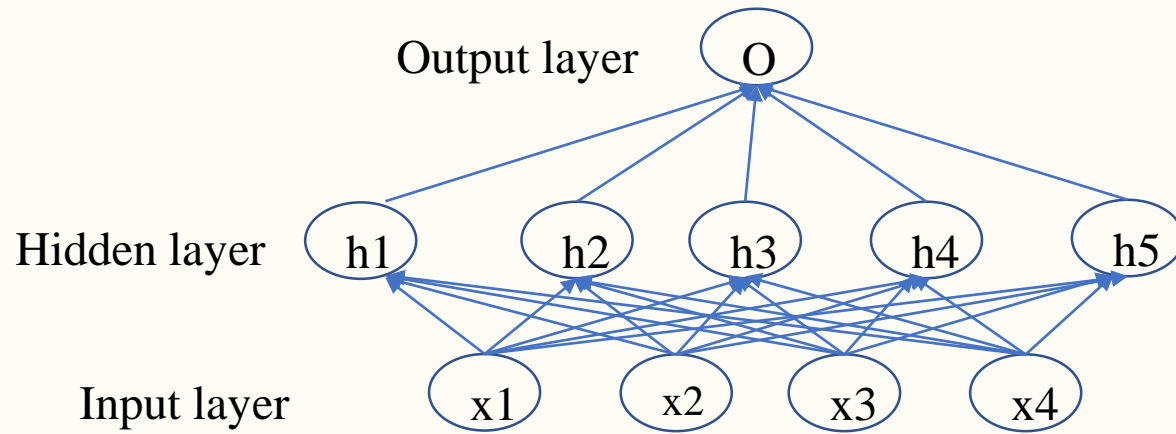
Logistic Regression (no bias)
1-layer model



2-layer model



Deep Neural Network



$$O = W_2 h \quad W_2 \text{ is } 1 \times 5$$

$$h = \sigma(W_1 x) \quad W_1 \text{ is } 5 \times 4$$

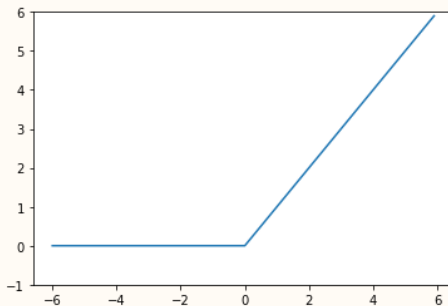
σ is a non-linear function
Applied elementwise

Use non-linear activation functions

Common activation functions

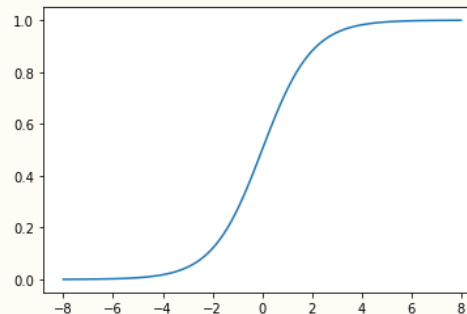
-Rectified Linear Unit (ReLU)

$$\text{ReLU}(z) = \max(z, 0)$$



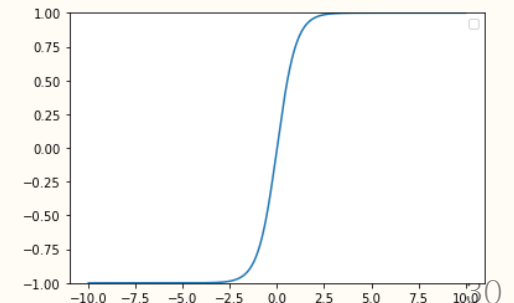
-Sigmoid

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

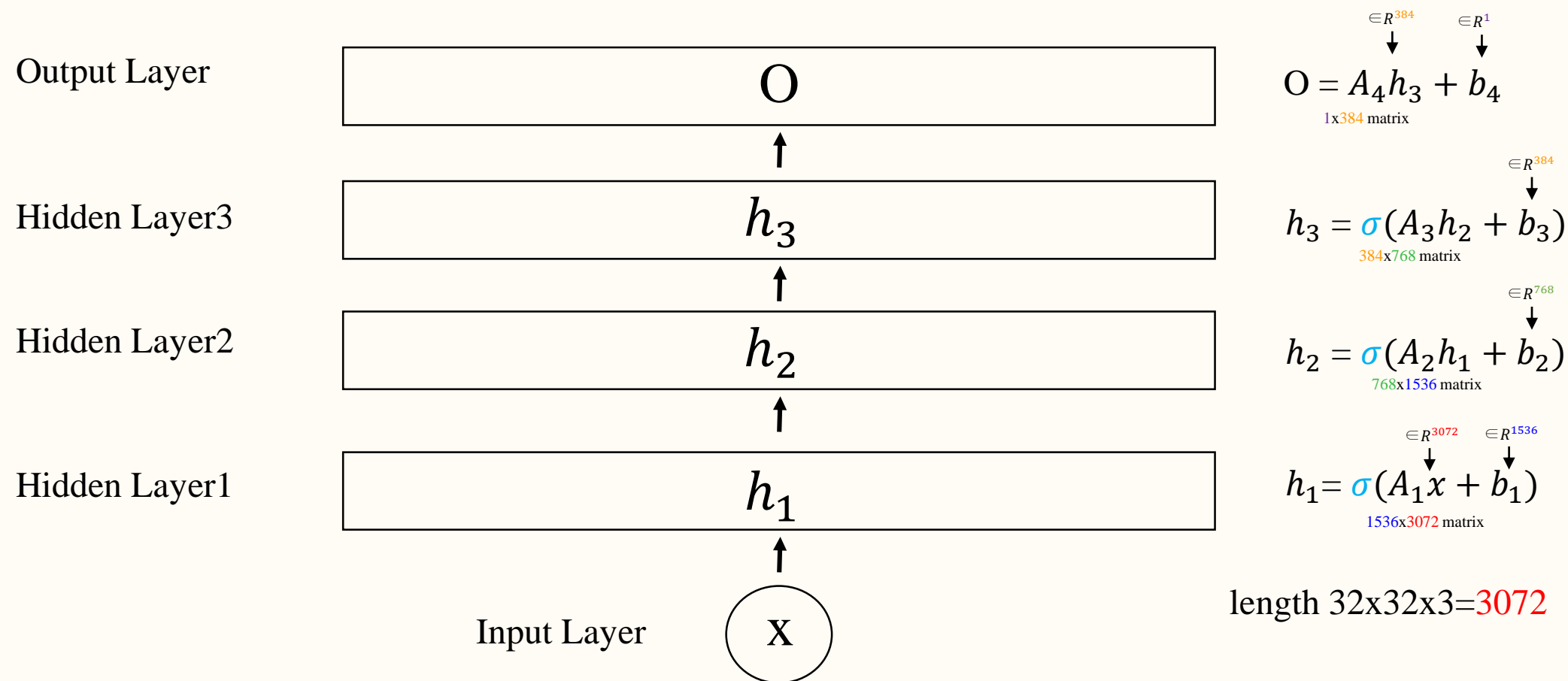


-Hyperbolic Tangent

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$



Architecture for CIFAR10 Binary classification



Activation function: $\sigma = \text{ReLU}$

Shift Invariance in Vision to Convolution

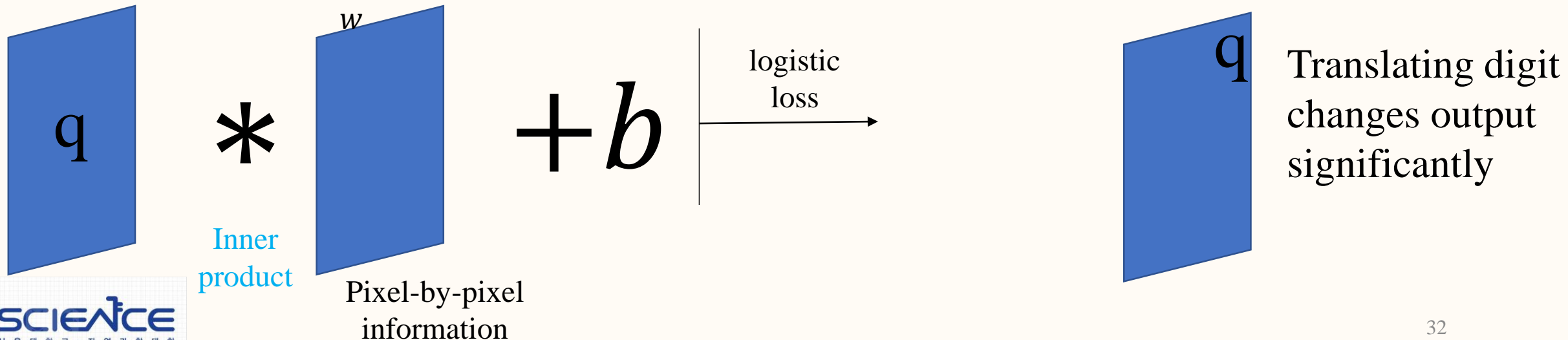


Cat



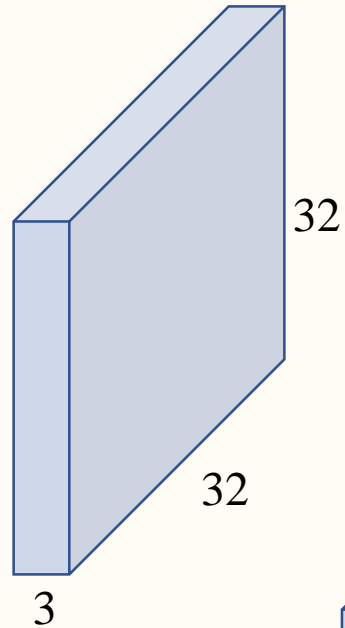
Still a Cat

Logistic regression (with a single fully connected layer) does not encode shift invariance

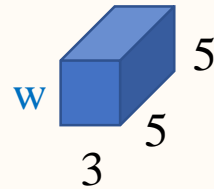


Convolutional Layer

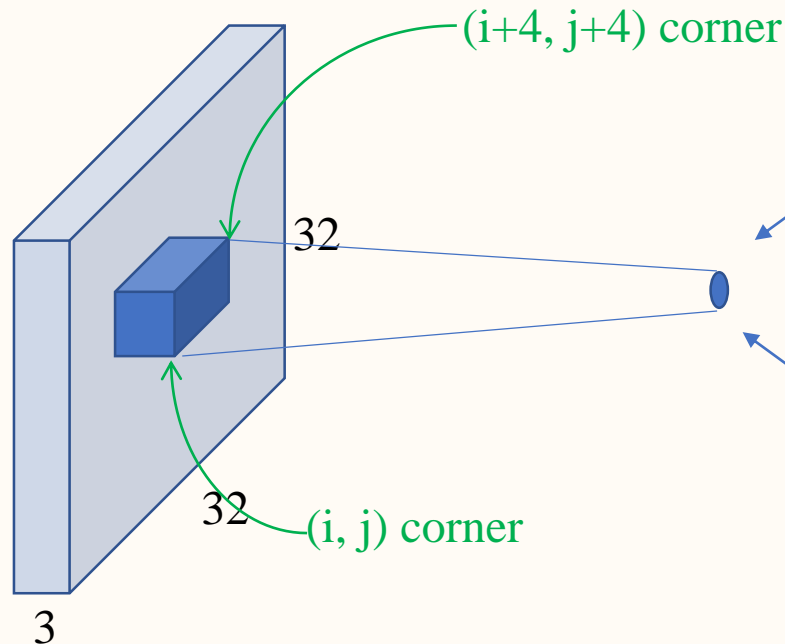
$3 \times 32 \times 32$ image



$3 \times 5 \times 5$ filter

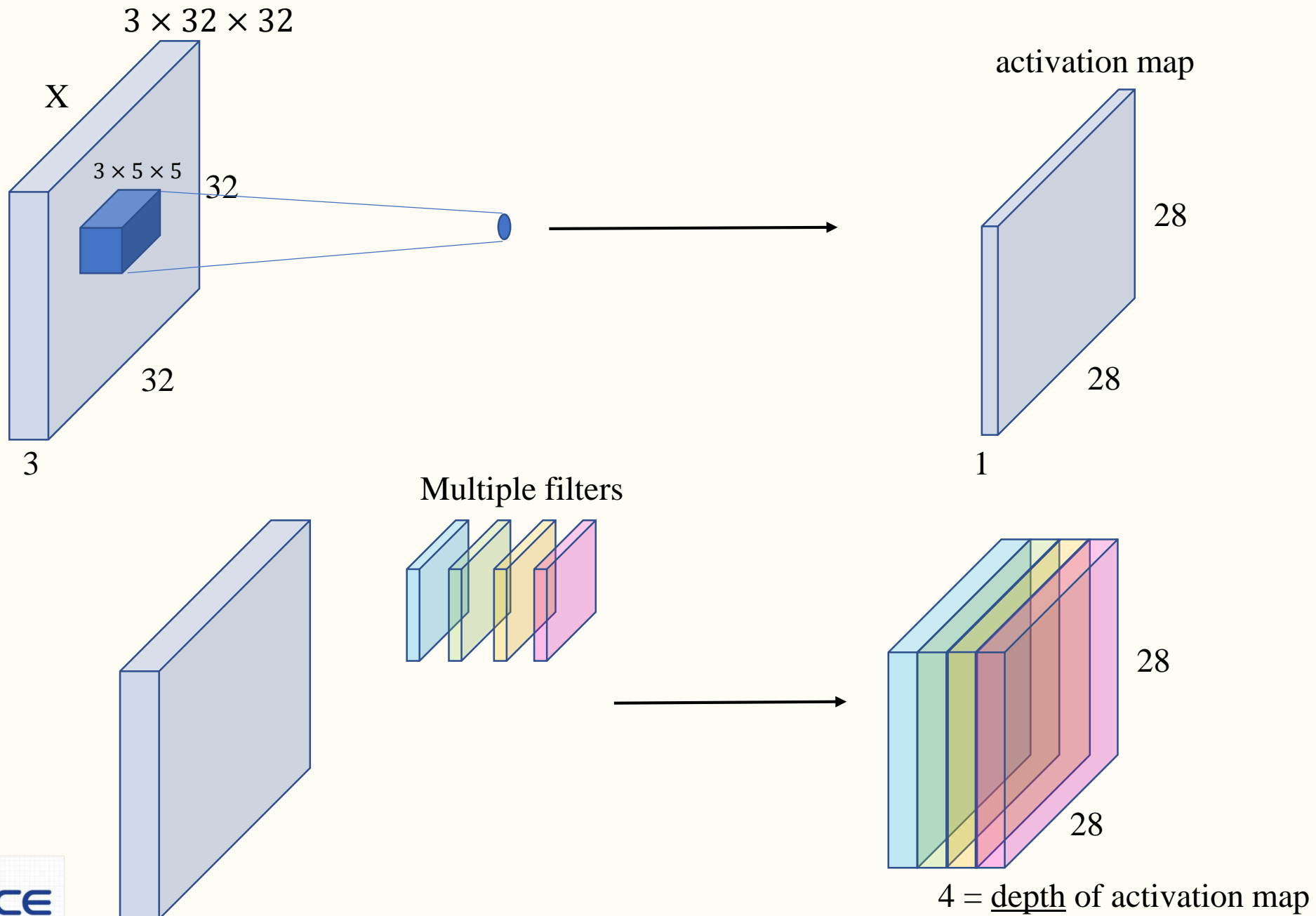


Convolve the filter with the image
(=Slide the filter spatially over the image and compute dot products)



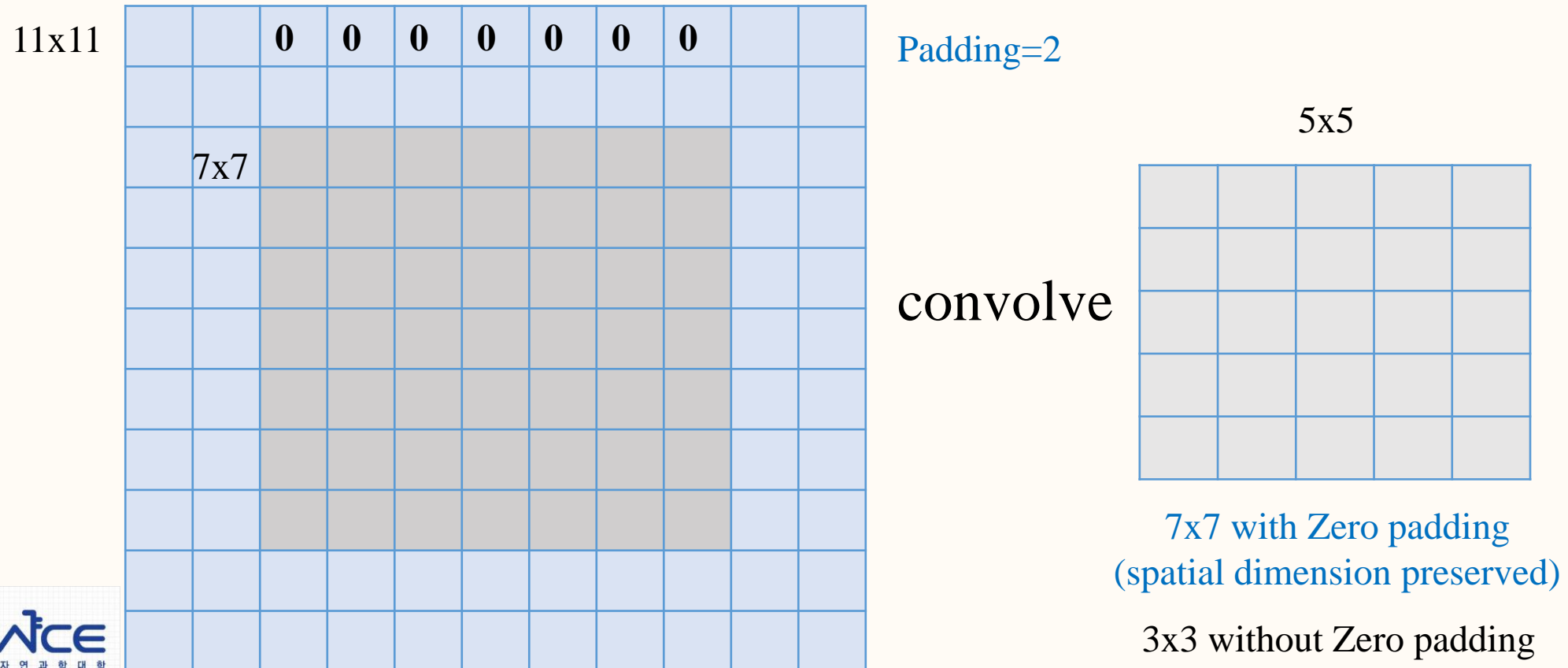
1 number: take a $3 \times 5 \times 5$ chunk of the image and take the inner product with w and add bias b

$= w.\text{reshape}(-1).\text{T} @ X[:, i:i+5, j:j+5].\text{reshape}(-1) + b$



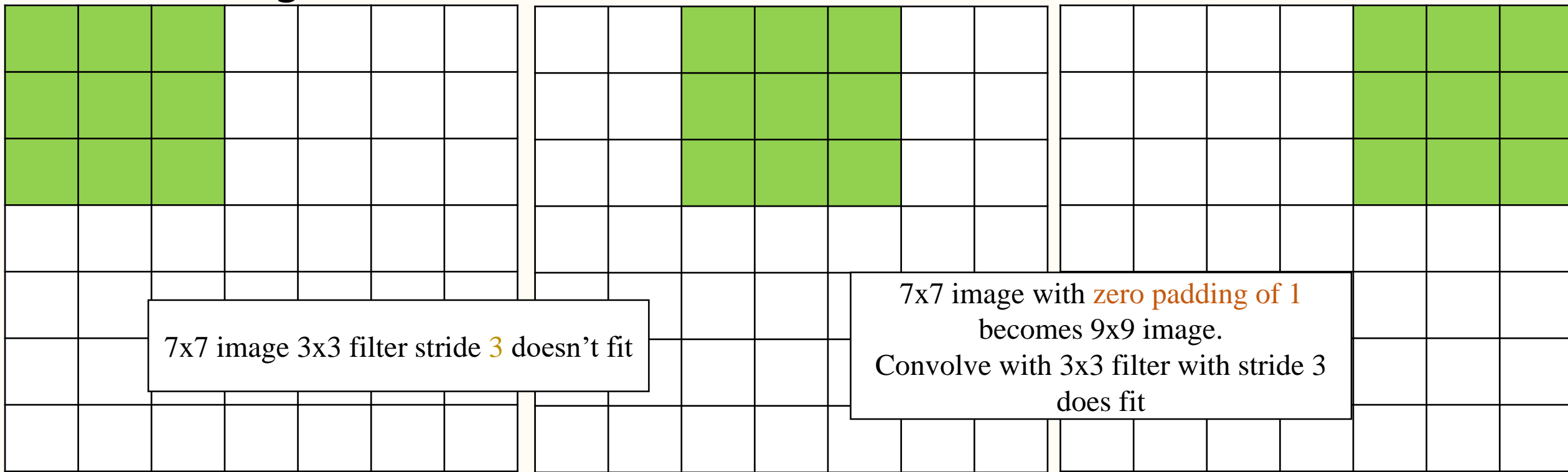
Convolution options : Zero Padding

$3 \times 32 \times 32$ convolved with $3 \times 5 \times 5$ filter $\Rightarrow 1 \times 28 \times 28$ activation map
Spatial dimension 32 reduced to 28



Convolution options : Stride

7x7 image convolved with 3x3 filter **with stride 2**



Output 3x3
(with stride 1, output is 5x5)

Summary

Input $W_1 \times H_1 \times D_1$

Conv layer parameters

\models K filters

\models F spatial extent ($F \times F \times D_1$ filters)

\models S stride

\models P padding

Output $W_2 \times H_2 \times D_2$

$$W_2 = \frac{W_1 - F + 2P}{S} + 1$$

$$H_2 = \frac{H_1 - F + 2P}{S} + 1$$

$$D_2 = K$$

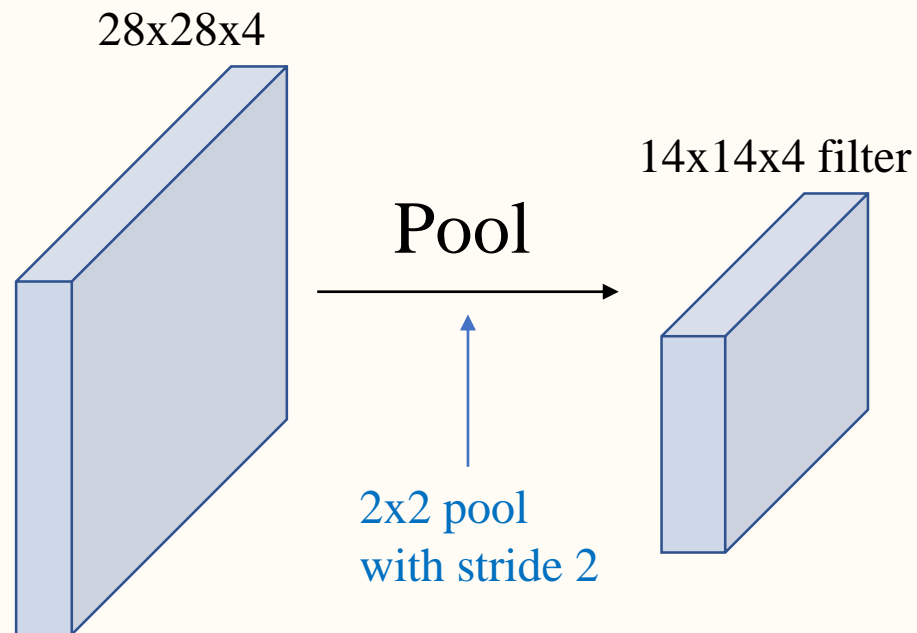
The number of parameters

$$F^2 D_1 K + K$$

filters biases

Pooling

- Similar to convolutions
- Used to reduce the size of the output
- Operates over each activation map independently



Single depth size

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

Max Pool

2x2 filters and stride 2

6	8
3	4

Not an instance
of convolution

Precise definitions in
`torch.nn.MaxPool2D`
`torch.nn.AvgPool2D`

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

Average Pool

2x2 filters and stride 2

Effect is subsampling
(lowering image resolution)

3.25	5.25
2	2

Instance of convolution
with fixed (untrainable)
weights, including the
independent operation
over each activation
map.
(Why?)

LeNet5

(LeCun, Bottou, Bengio, Haffner 1998)

Modern instances of LeNet5 use

- $\sigma = \text{ReLU}$
- MaxPool instead of avg. pool
- No σ after S2, S4 (Why?)
- No Gaussian connections
- Complete C4 connections

28x28 MNIST image
with $p=2 \Rightarrow 32 \times 32$

