**MACHINE LEARNING DAY 2**

# DEEP LEARNING

## Session IV: Convolutional Neural Network

**Isaac Ye, HPTC @ York University**

**Isaac@sharcnet.ca**

# Session IV

- Backward propagation

- Model capacity / overfitting

- MNIST classification

- Vanishing gradient problem

- *Lab 4A: Multi-Layer Perceptron (MNIST)*

- Issue with MLP

- Convolutional Neural Network

- Techniques

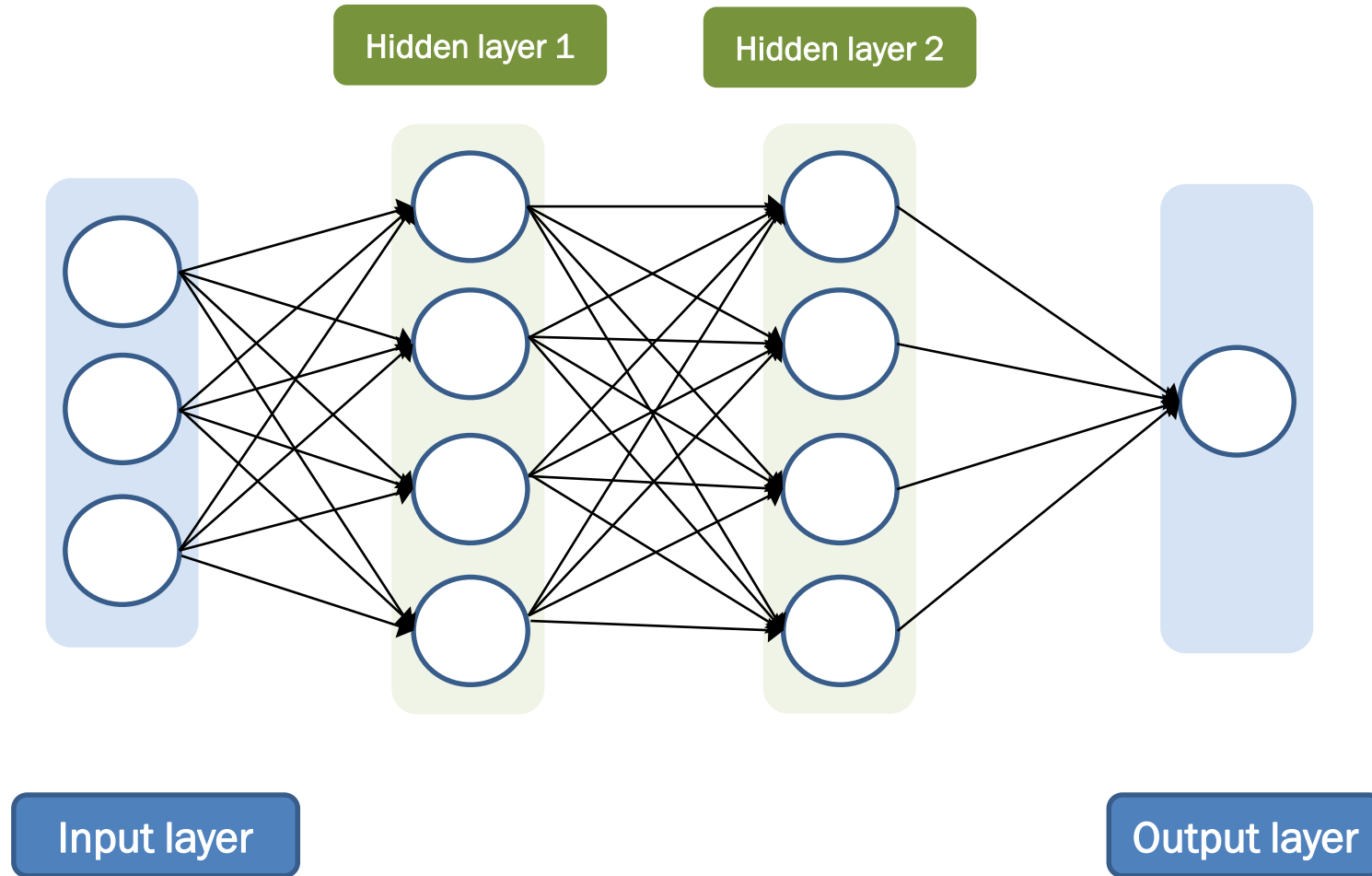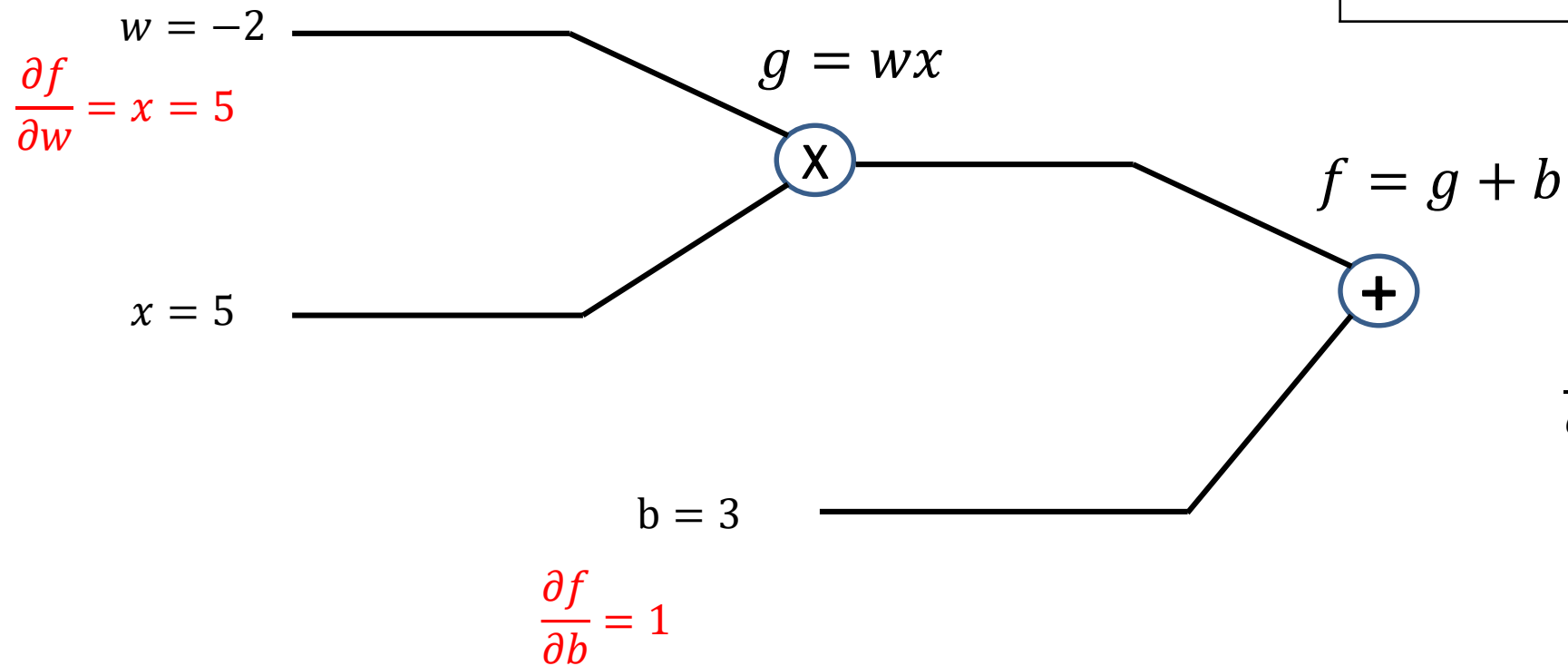- *Lab 4B: Convolutional Neural Network (CIFAR10)*

# Supervised learning

| | Supervised |
|---|---|
| **Discrete** | Classification |
| **Continuous** | Regression |

## Supervised

| | |
|---|---|
| Input | Data (x, y)    x: data, y: label |
| Goal | learn a function to map x to y |
| Examples | classification, regression, object detection |
| | semantic segmentation, image captioning |

# Categories of ML problems

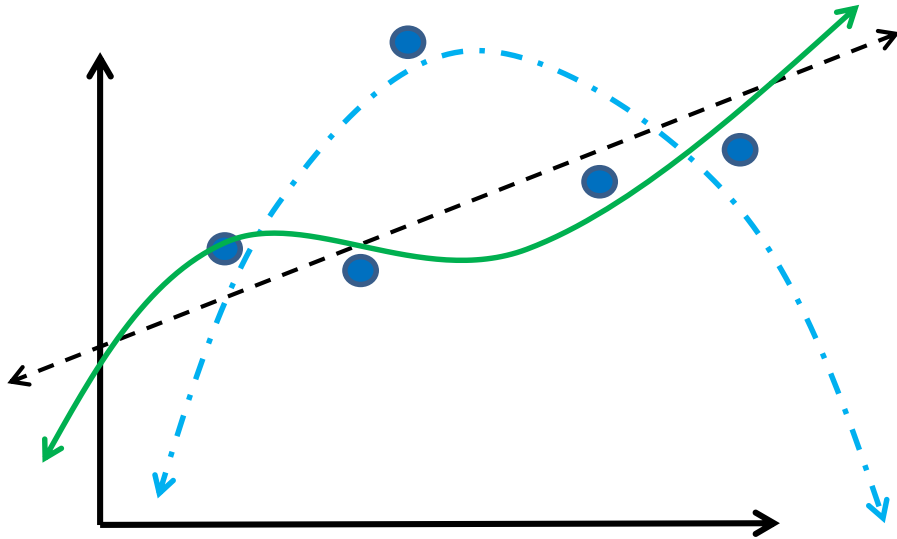|  | Supervised | Unsupervised | Reinforcement |
|---|---|---|---|
| **Discrete** | Classification | Clustering | Action space agent |
| **Continuous** | Regression | Dimensionality reduction | Action space agent |

# Multi-Layer Perceptron



Hidden layer 1

Hidden layer 2

Input layer

Output layer

# Backward propagation

| $g = wx$ | $f = g + b$ |
|---|---|
| $\frac{\partial g}{\partial w} = x \,, \frac{\partial g}{\partial x} = w$ | $\frac{\partial f}{\partial g} = 1 \,, \frac{\partial f}{\partial b} = 1$ |

$w = -2$

$\frac{\partial f}{\partial w} = x = 5$

$g = wx$

(X)

$x = 5$

$f = g + b$

(+)

$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w} = 1x$

$b = 3$

$\frac{\partial f}{\partial b} = 1$

# Model capacity

$$y = w_1 x + b$$

$$y = w_2 x^2 + w_1 x + b$$

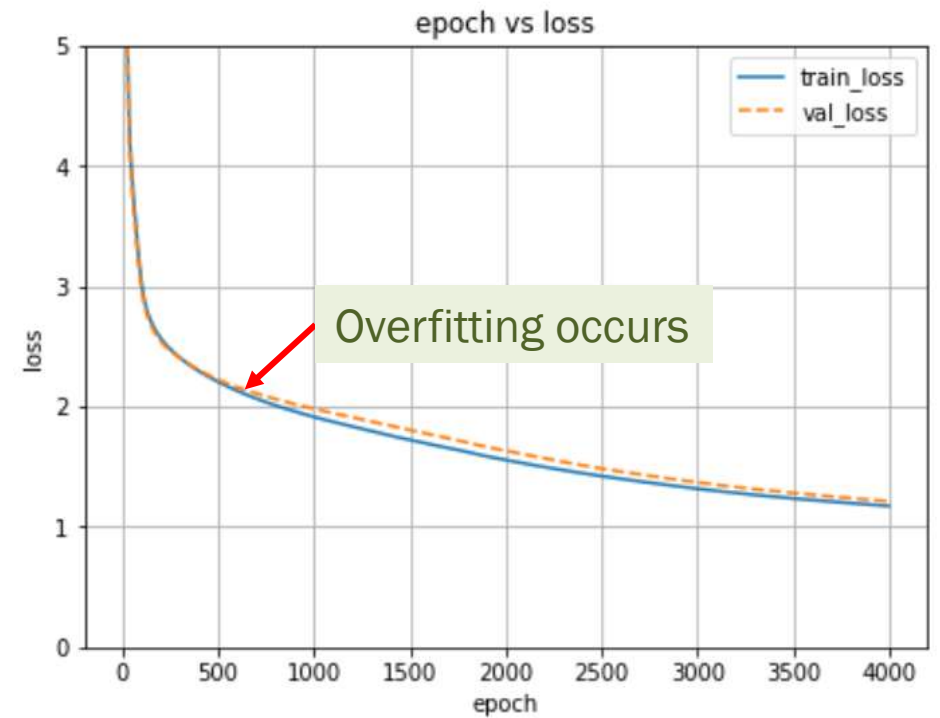$$y = w_3 x^3 + w_2 x^2 + w_1 x + b$$

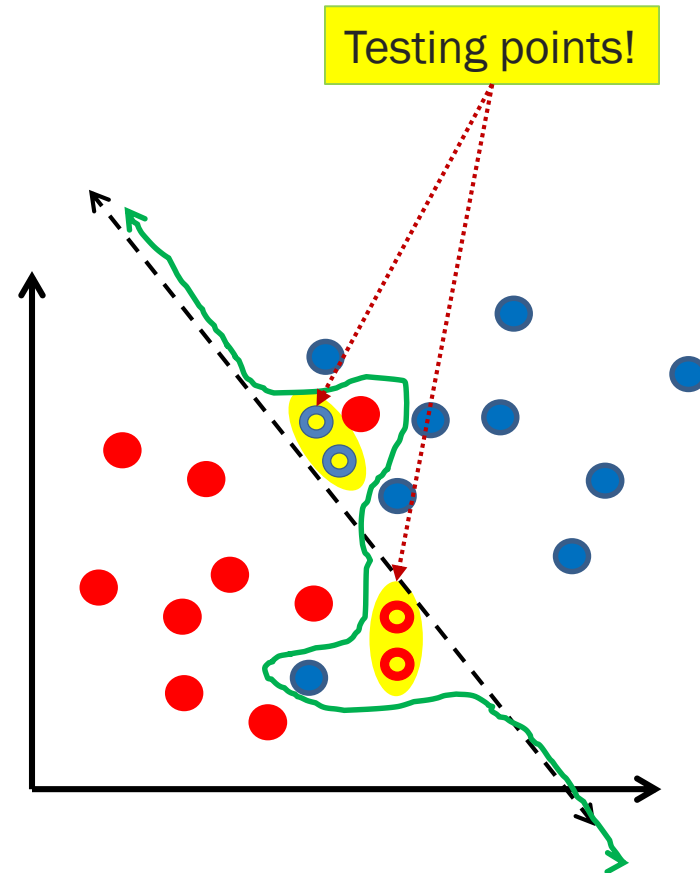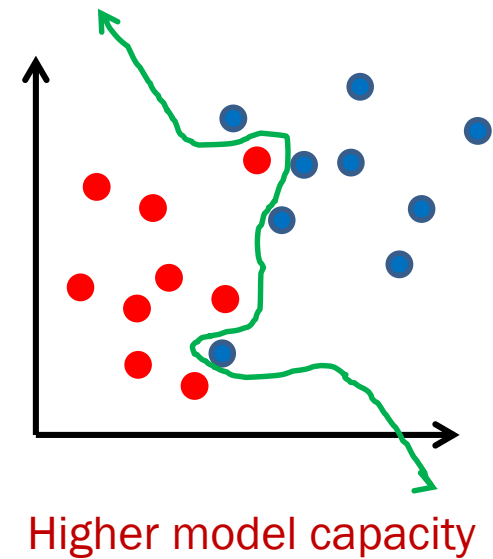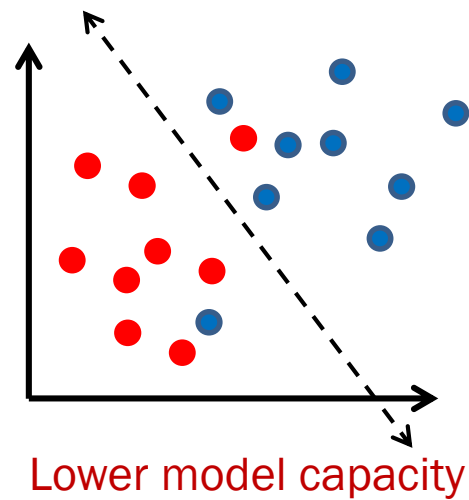Higher model capacity

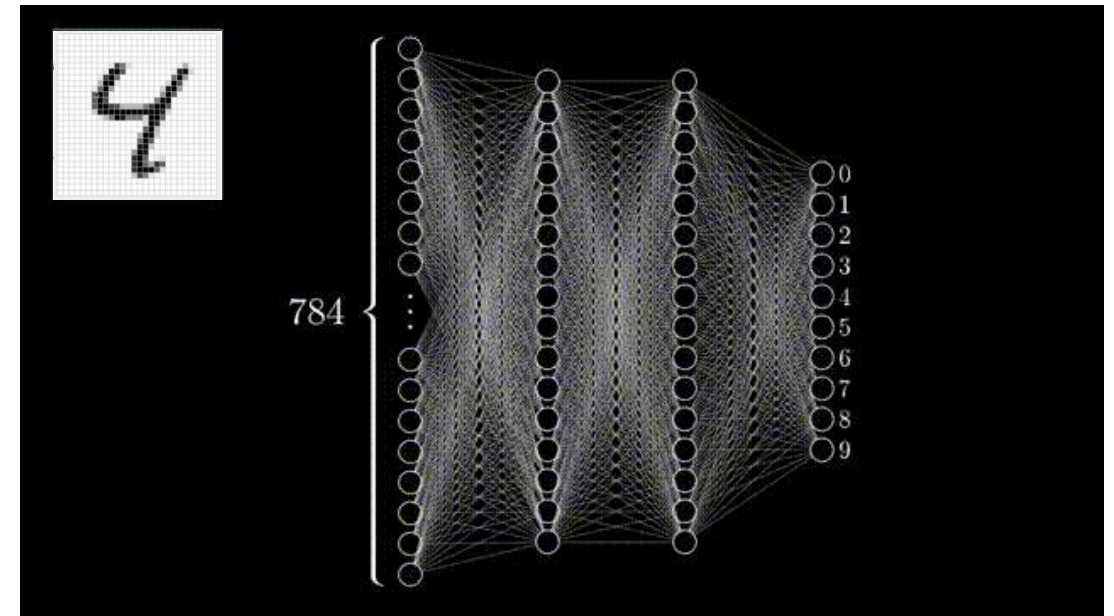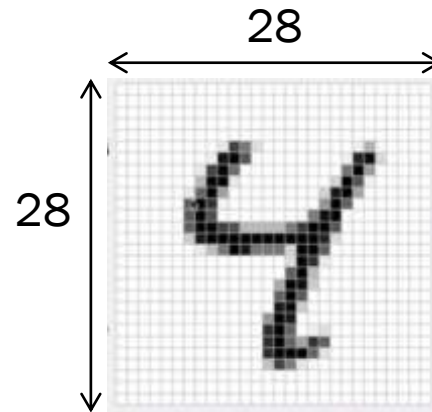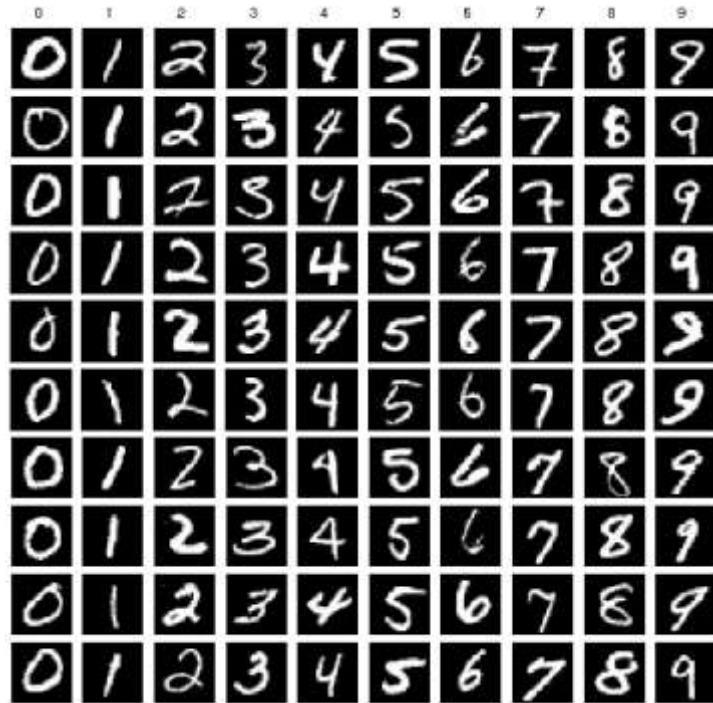The more hidden layers and units,
the higher model capacity it will have.

Does it guarantee better accuracy?

Hidden layer 1     Hidden layer 2

Input layer     Output layer

# Overfitting



Lower model capacity

Higher model capacity

Testing points!

epoch vs loss

Overfitting occurs

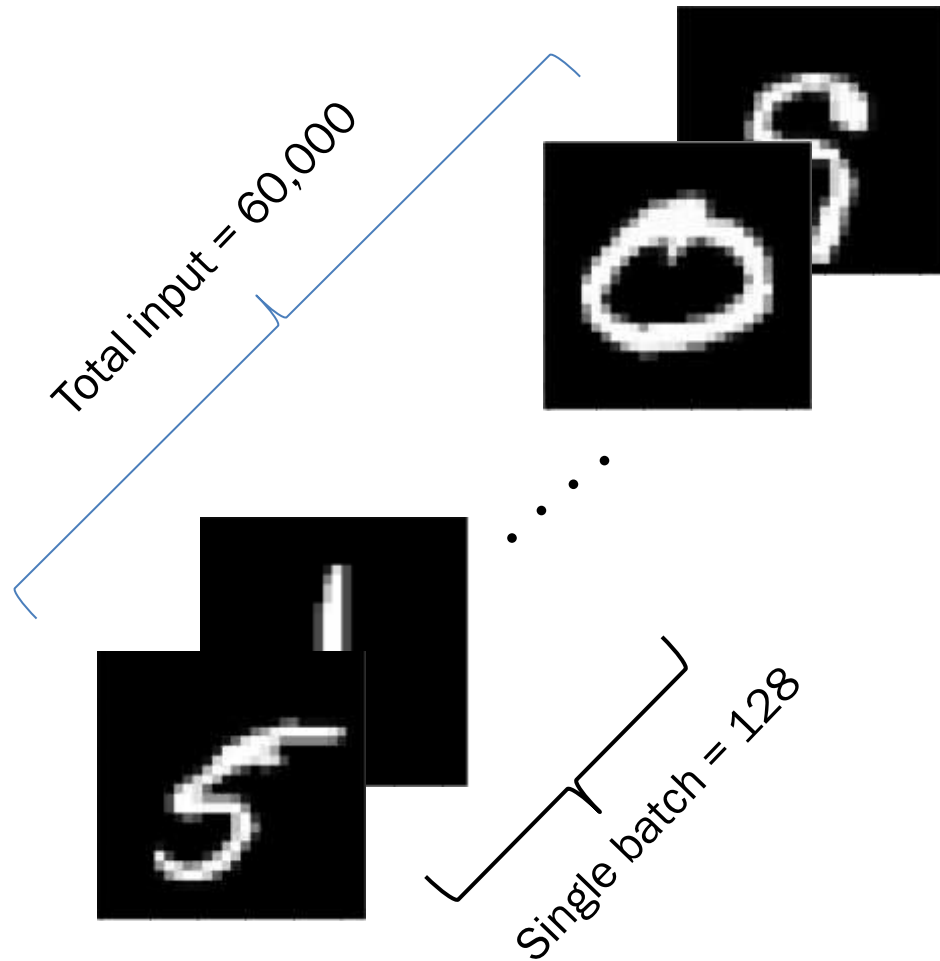# Classification problem: MNIST



28

28

784

0
1
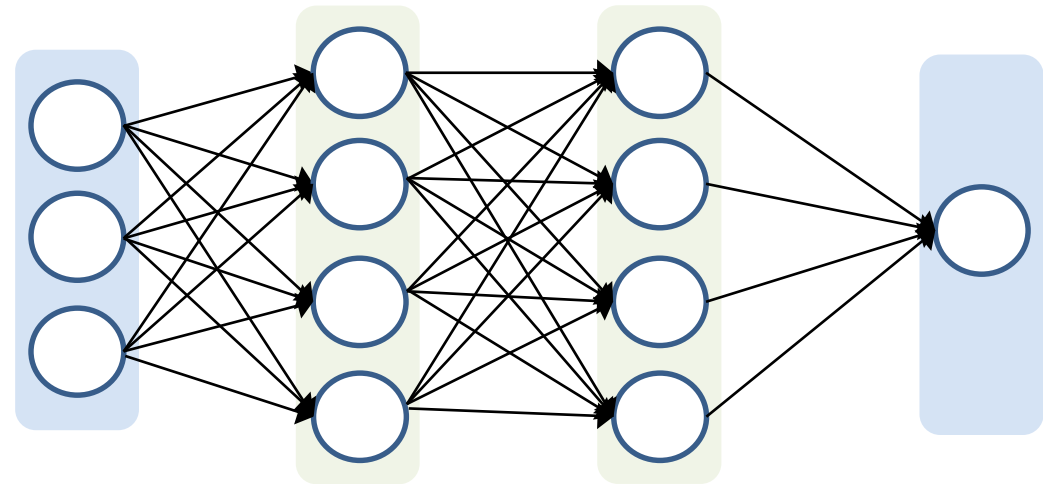2
3
4
5
6
7
8
9

Handwritten data
60K train set and 10K test set
Each image has a size of 28x28 (=784)

# Batch, Iteration, Epoch



Total input = 60,000

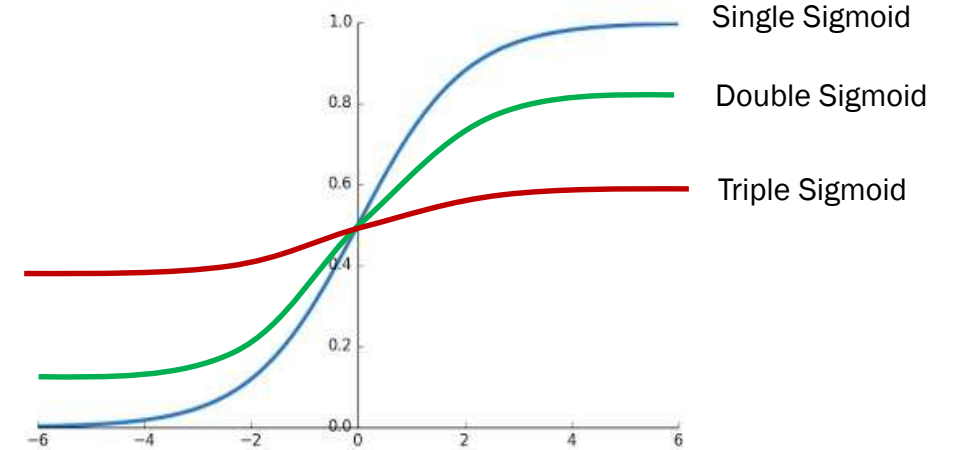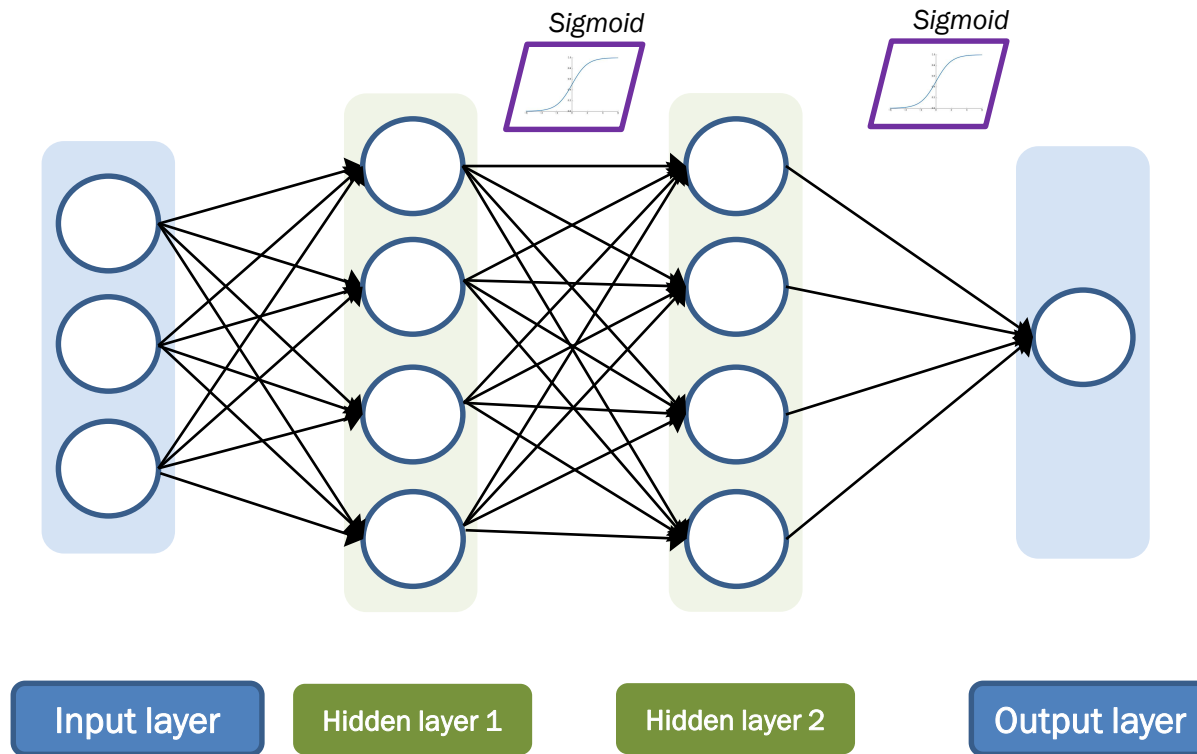Single batch = 128
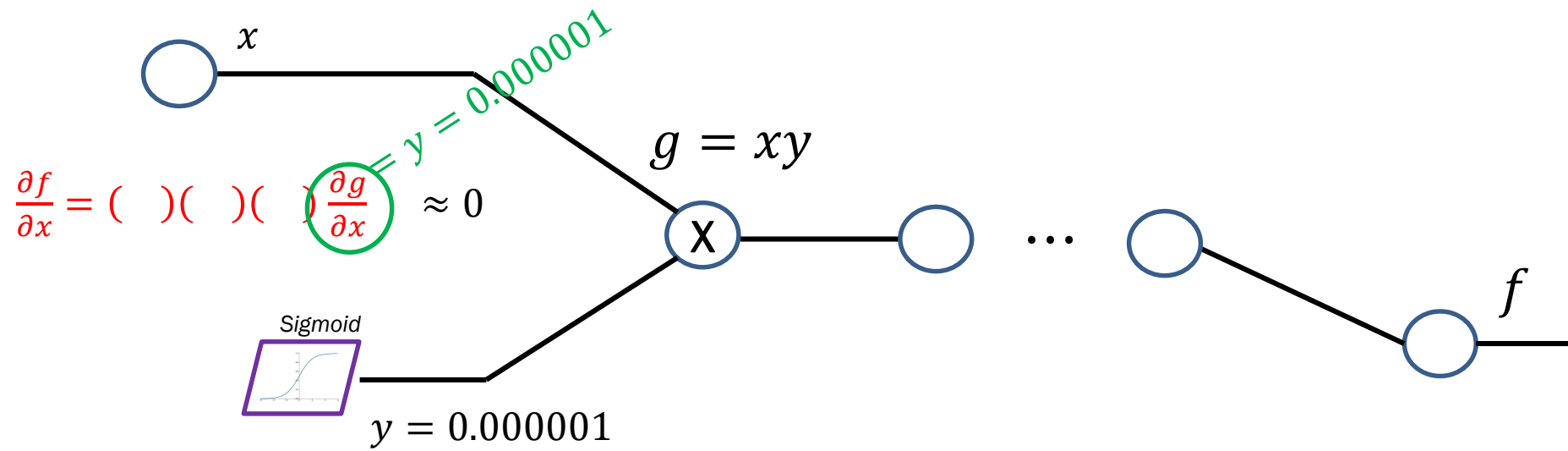
Epoch: one complete run of total input
Batch size: the amount of input for each iteration
# of iteration = # total input / batch size

# Activation function: Sigmoid problem
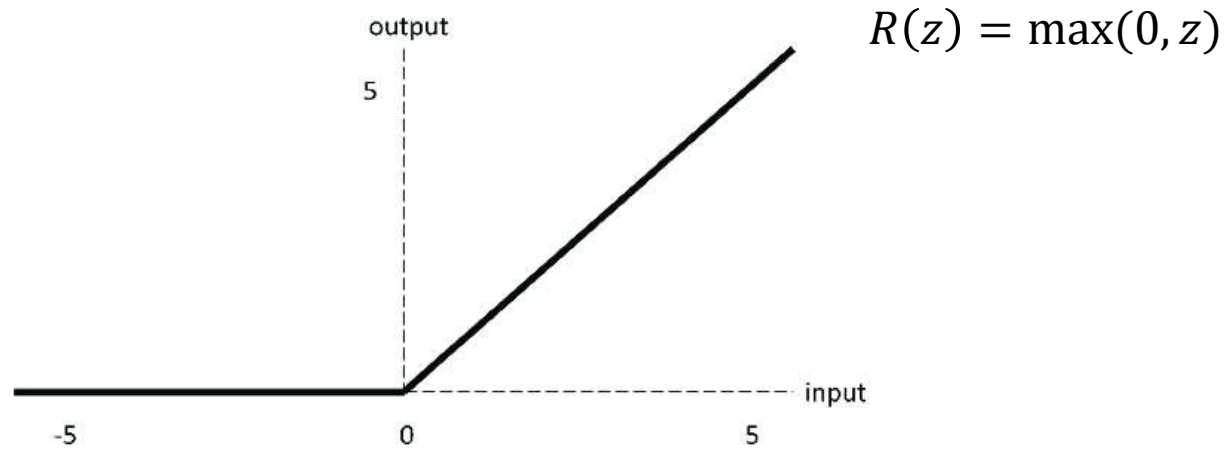


| Input layer | Hidden layer 1 | Hidden layer 2 | Output layer |

# Vanishing gradient problem

# Activation functions: Rectified Linear Unit (ReLU)



$$R(z) = \max(0, z)$$

# Lab 4A: MNIST classification - MLP

| Exercise 1: MLP vs. Linear model | Model |
|---|---|
| Run 1 | MLP |
| Run 2 | Linear model |

| Exercise 2: more layers w/ 100 units | # of layers |
|---|---|
| Run 1 | 2 |
| Run 2 | 4 |

| Exercise 3: Different units | # of units |
|---|---|
| Run 1 | 100 |
| Run 2 | 200 |

| Exercise 4: different learning rate | Learning rate |
|---|---|
| Run 1 | 0.002 |
| Run 2 | 0.02 |

https://pytorch.org/docs/stable/nn.html#non-linear-activations-weighted-sum-nonlinearity
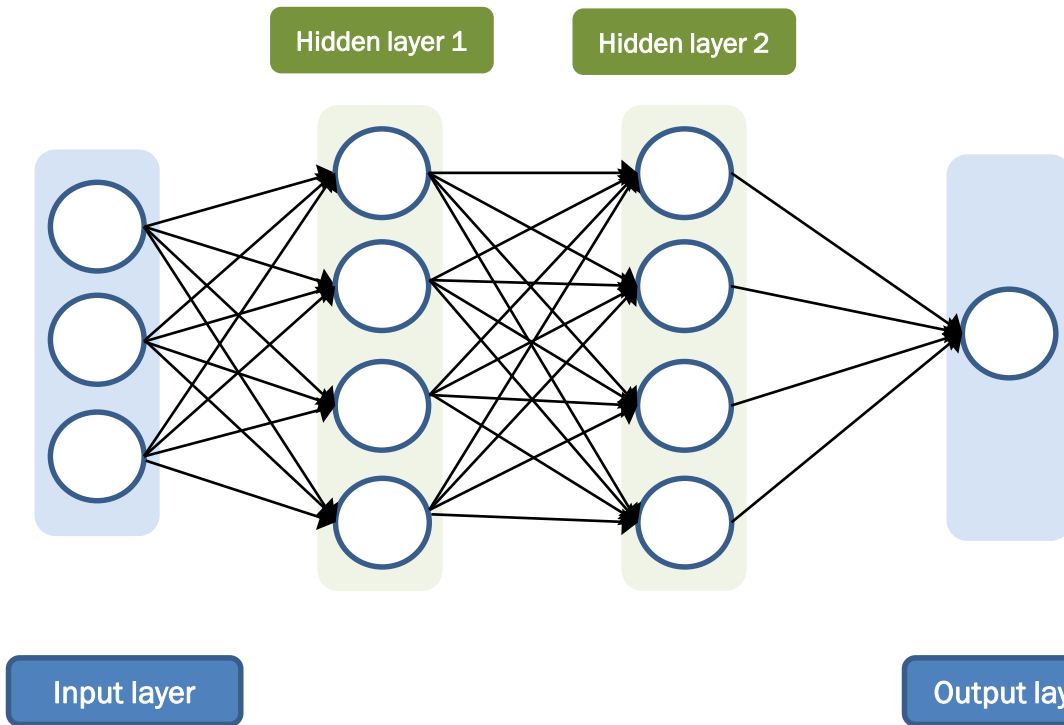
**You may want to try**

1. Use different cost function

2. Use different optimizers

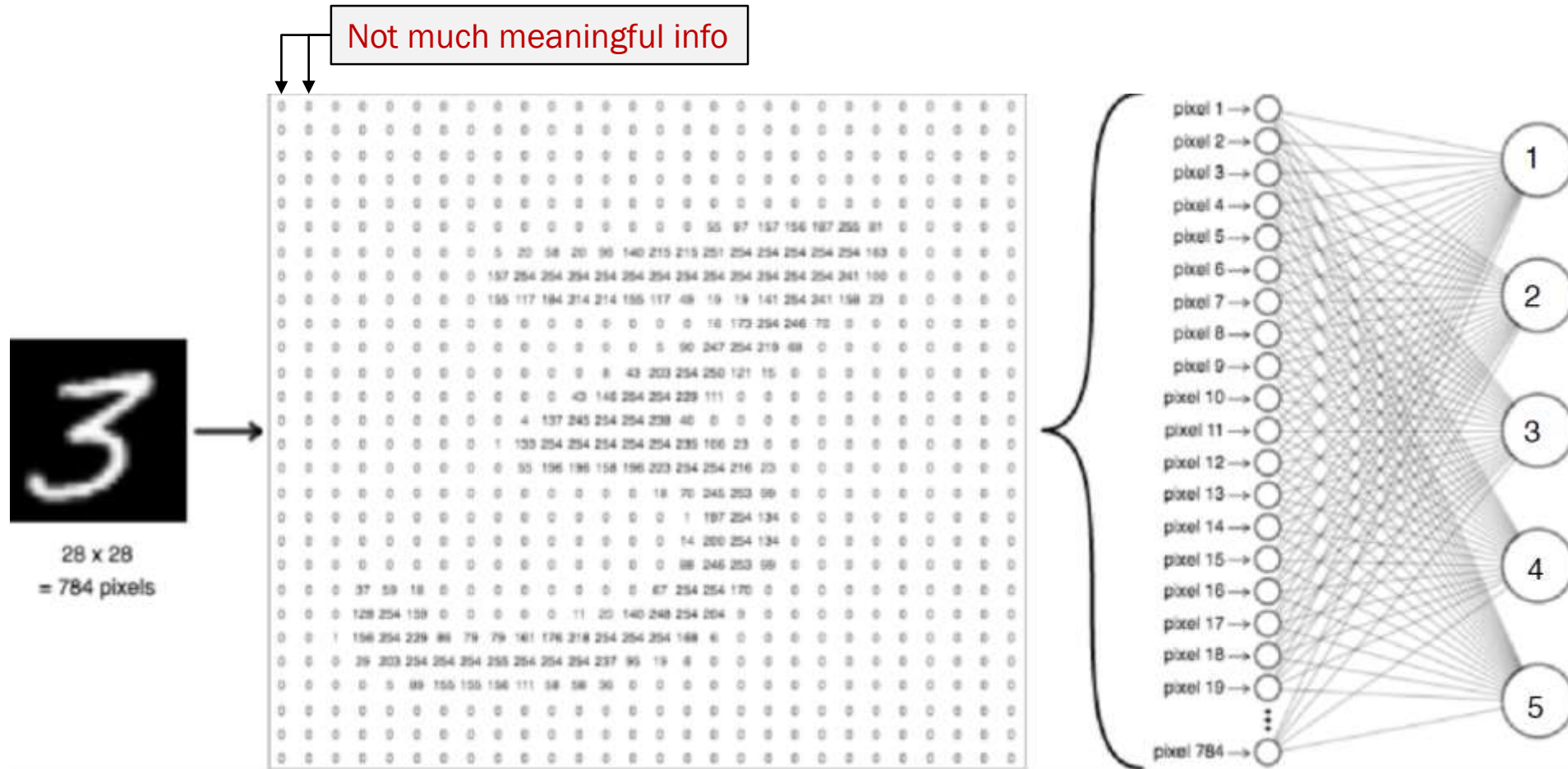   https://pytorch.org/docs/stable/optim.html?highlight=optimizer#torch.optim.Optimizer

Break room

# Issue with MLP



Hidden layer 1

Hidden layer 2

Input layer

Output layer

A neuron is connected with every neuron in next layer (fully connected)

# of parameters increases explosively

# Issue with MLP



Not much meaningful info

28 x 28
= 784 pixels

Some of parameters are meaningless!

# How to recognize an image?

# Hierarchical structure

Retinal ganglion cell receptive fields

LGN and V1 simple cells
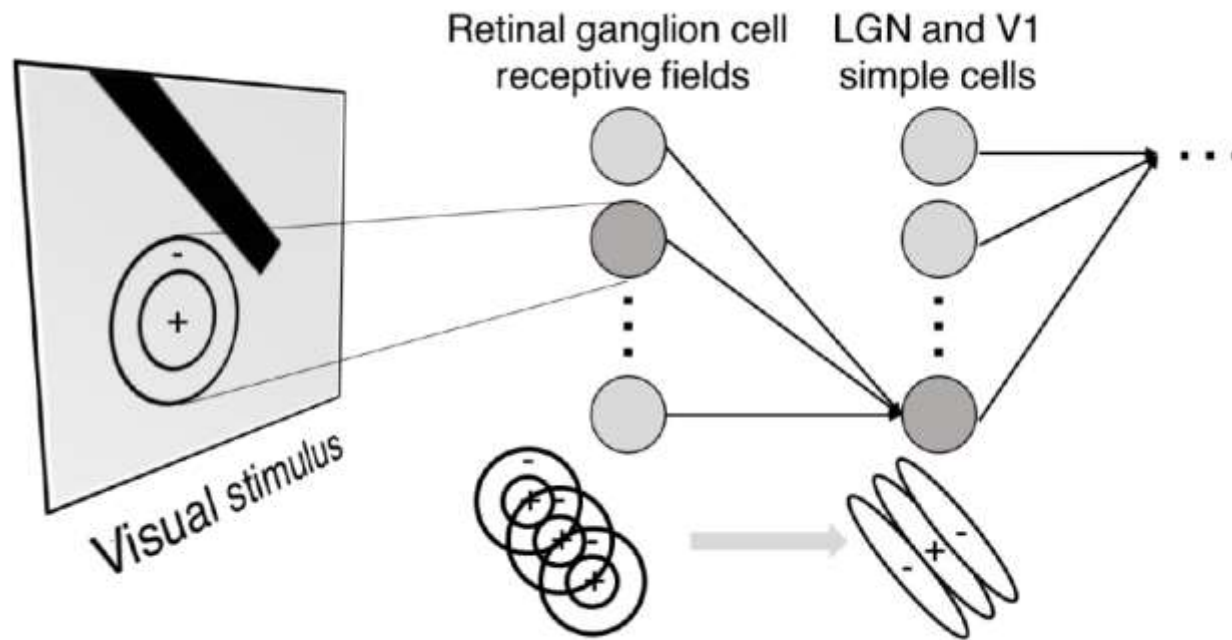
Visual stimulus

Illustration of hierarchical organization in early visual pathways by Lane McIntosh, copyright CS231n 2017

<u>Simple cells</u>: Response to light orientation

<u>Complex cells</u>: Response to light orientation and movement

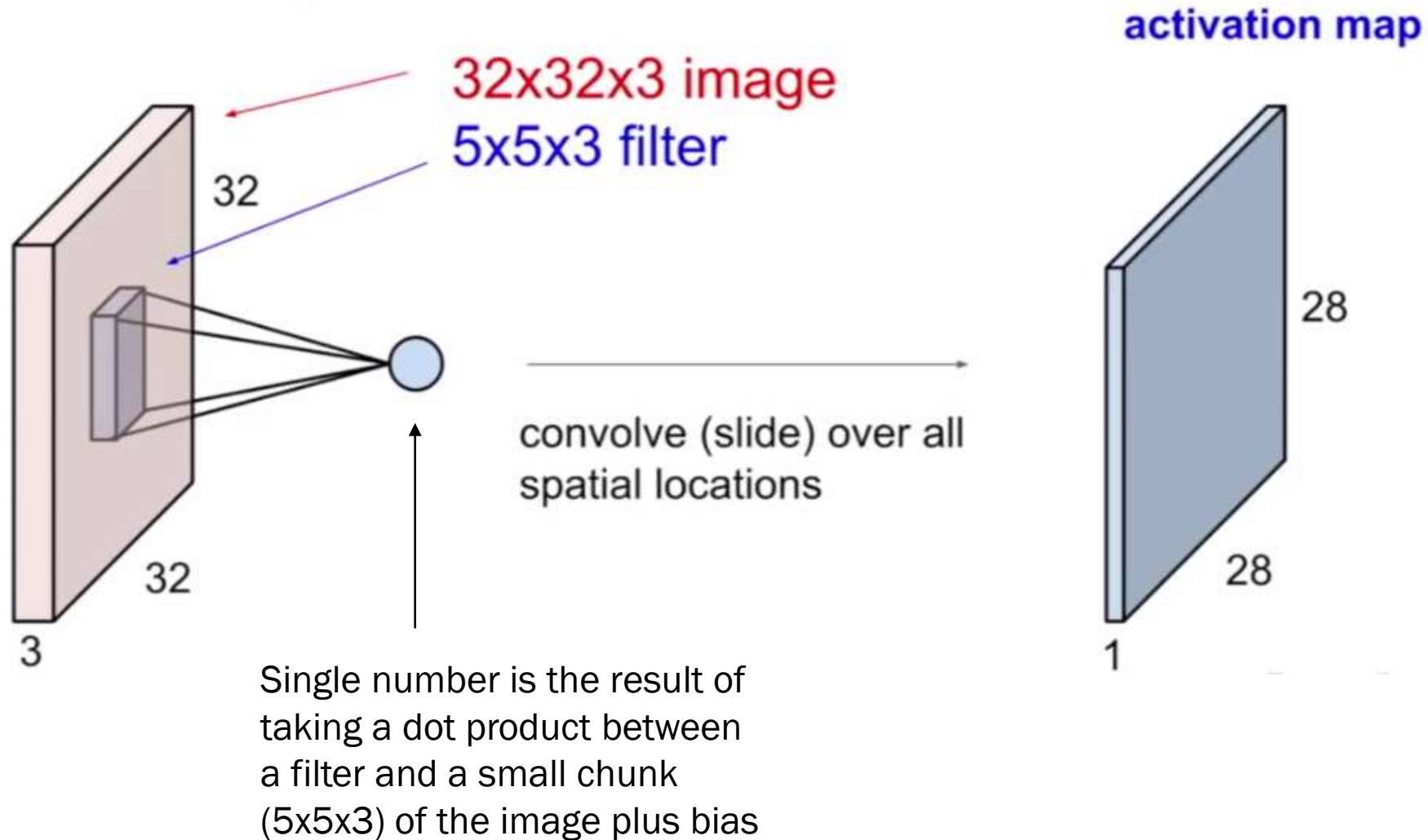<u>Hypercomplex cells</u>: Response to movement with an end point
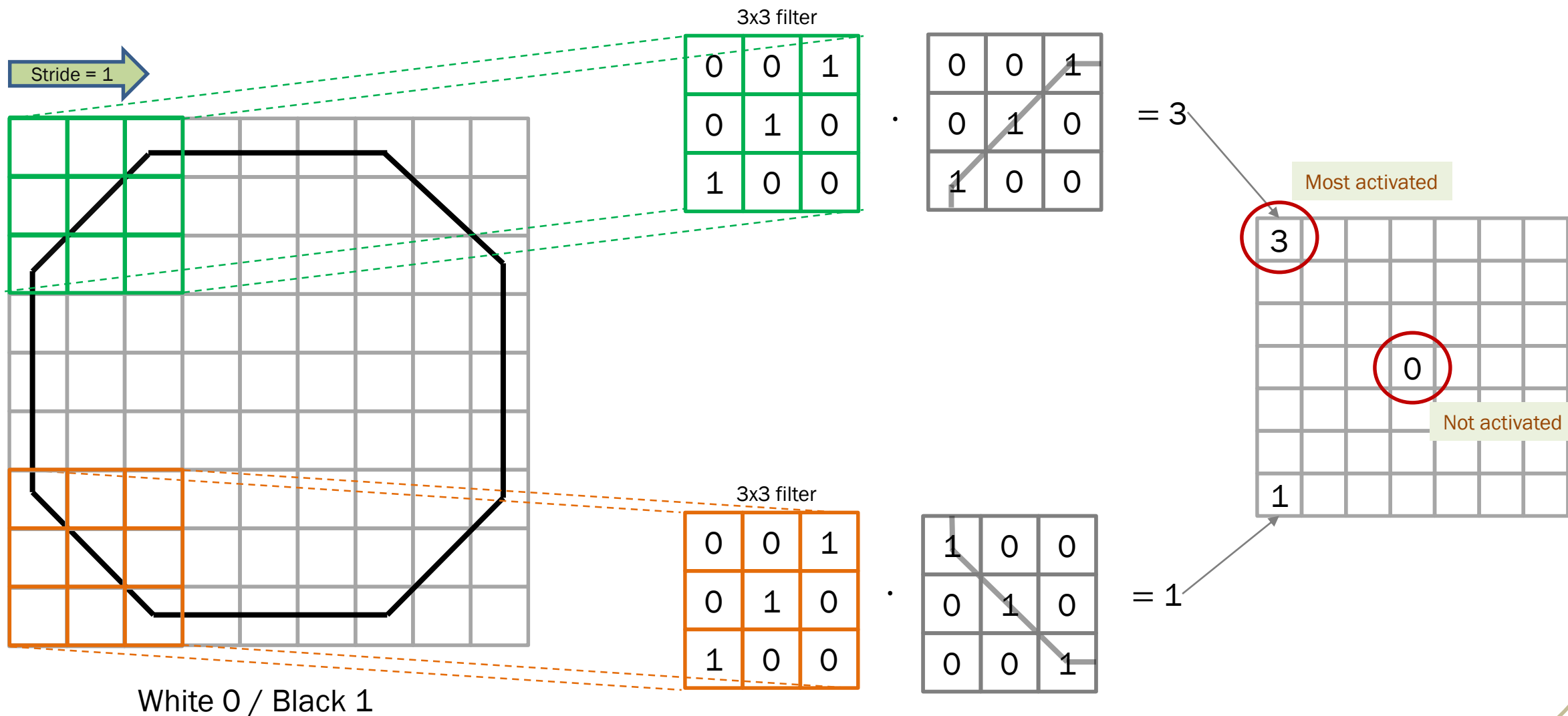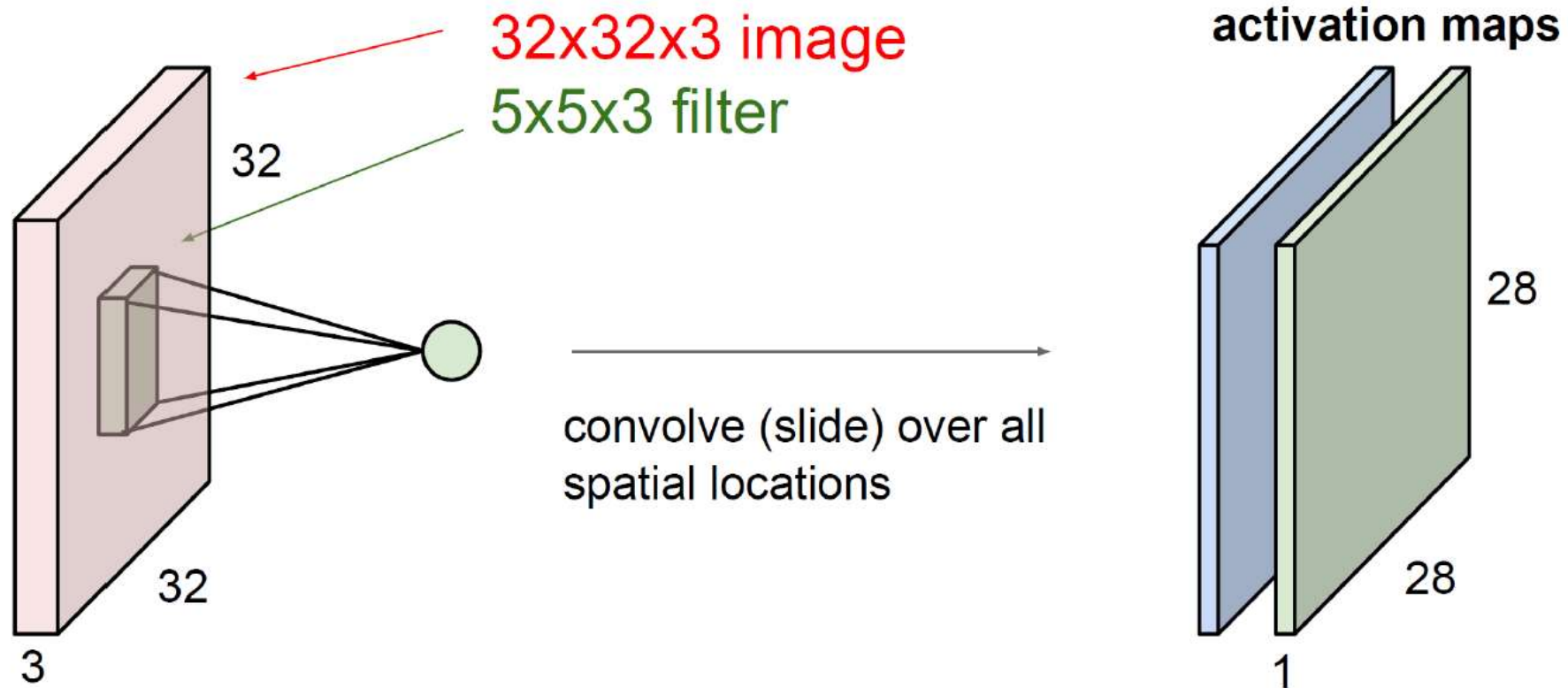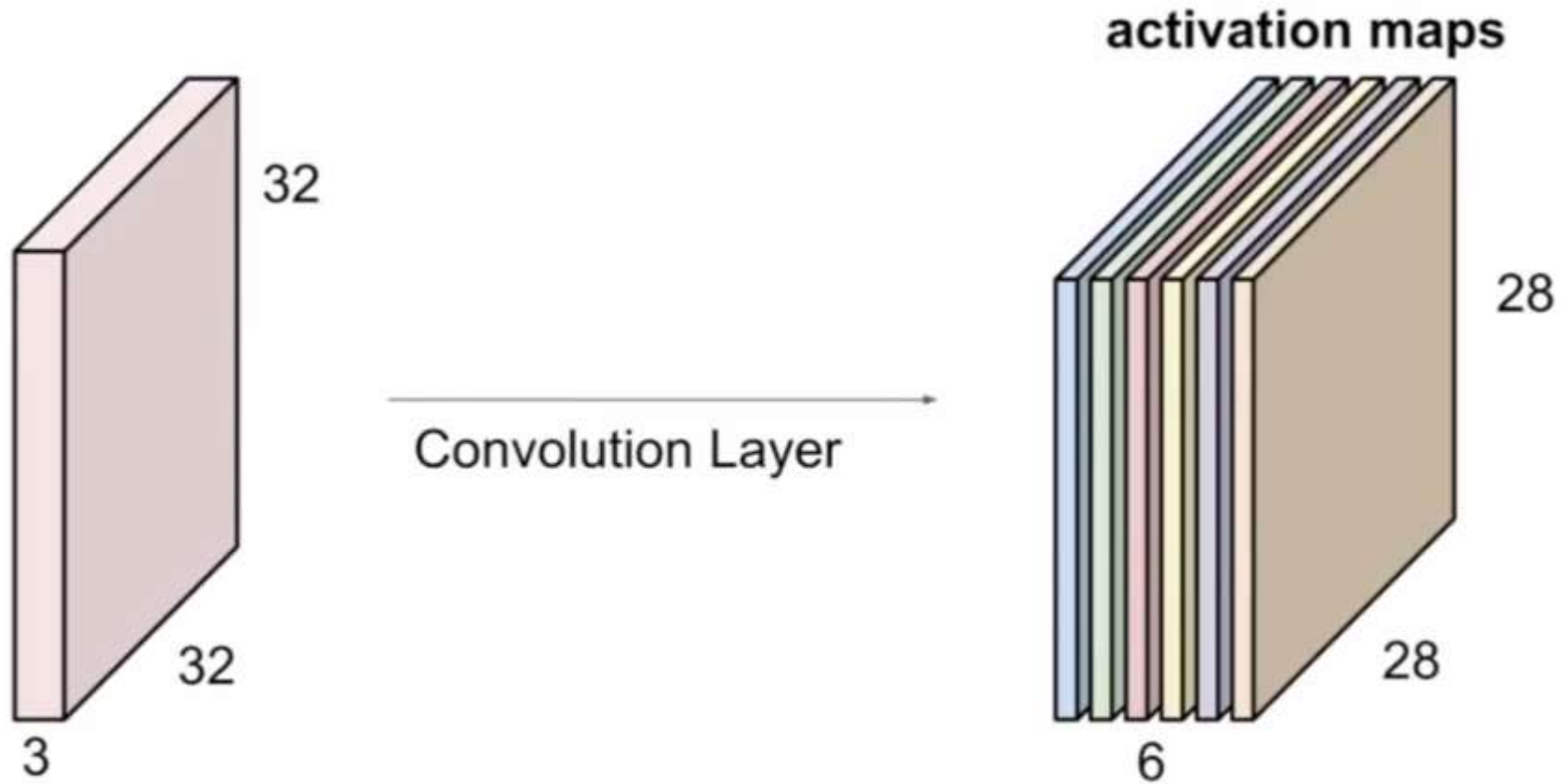
No response

Response (end point)

# Convolutional layer



32x32x3 image
5x5x3 filter

32

32

3

activation map

28

28

1

convolve (slide) over all spatial locations

Single number is the result of taking a dot product between a filter and a small chunk (5x5x3) of the image plus bias

# Convolutional operation

# Convolutional layer



32x32x3 image
5x5x3 filter

activation maps

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Convolutional layer



activation maps

32
32
3

Convolution Layer

28
28
6

# Convolutional Neural Network

# Pooling



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

4x4

**Max pooling**
filter: 2x2
Stride: 2

→

|   |   |
|---|---|
| 6 | 8 |
| 3 | 7 |

2x2

# ConvNet Architecture

# Hyperparameter?

- Non-learnable parameters

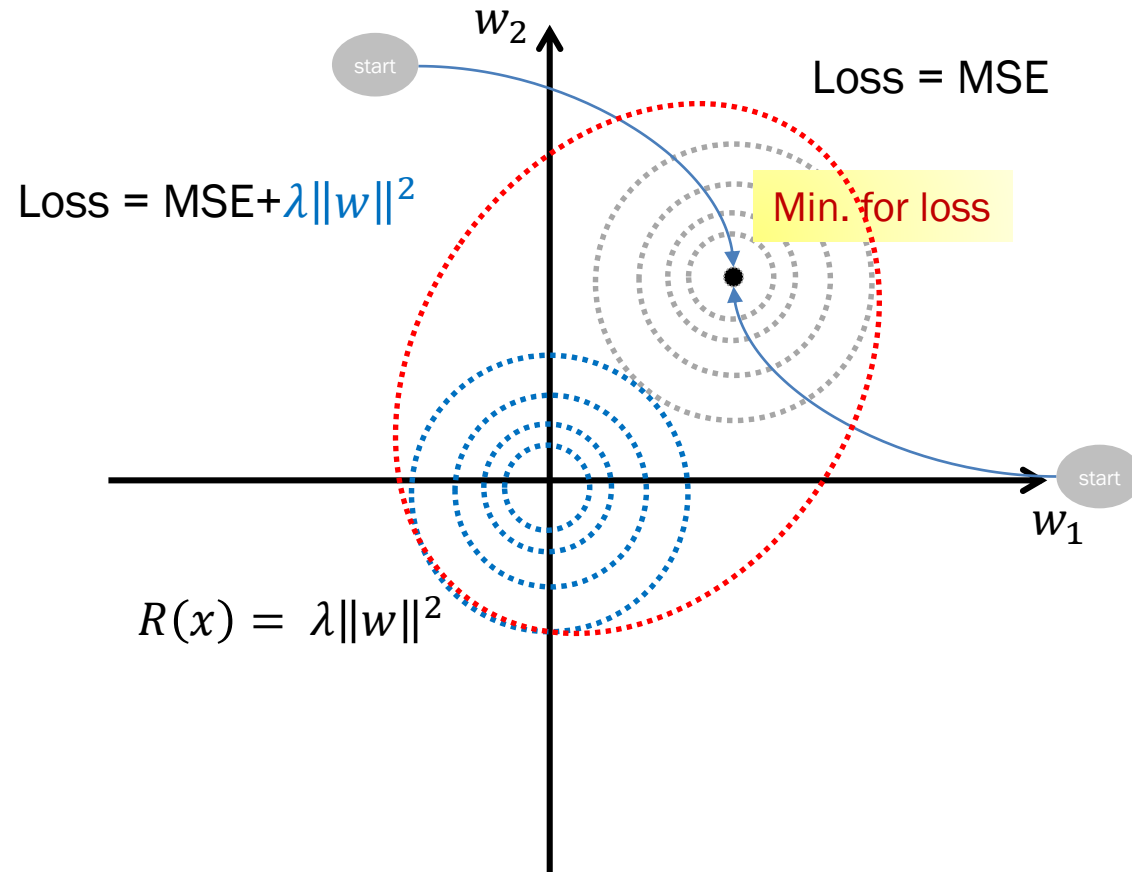| Category | Parameters |
|----------|------------|
| Model capacity | # of hidden layers<br># of hidden units<br>Activation function |
| Regularization | Dropout rate<br>Batch normalization<br>L2 regularization<br>Xavier initialization |
| Optimizing | Optimizer<br>Learning rate<br># of Epoch<br>Batch_size |
| Device | CPU/GPU |
| Post processing | Saving/filename |

# Dropout



Intentionally turn off nodes with probability when training
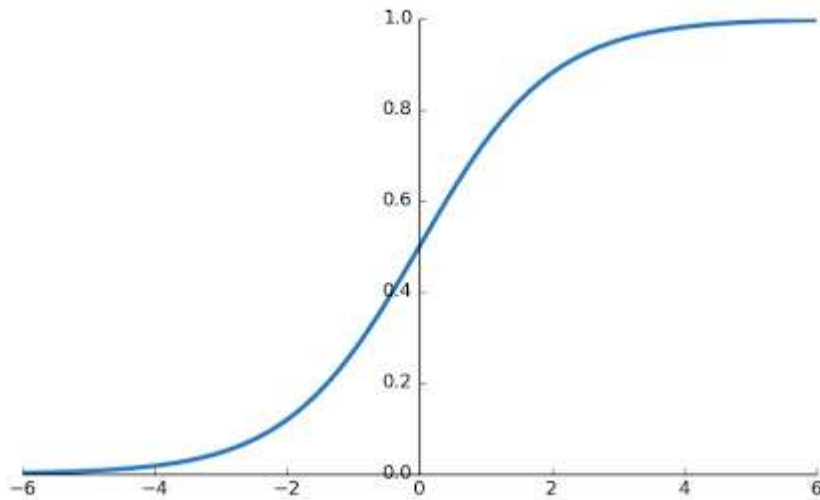
# L2 Regularization

# Xavier initialization
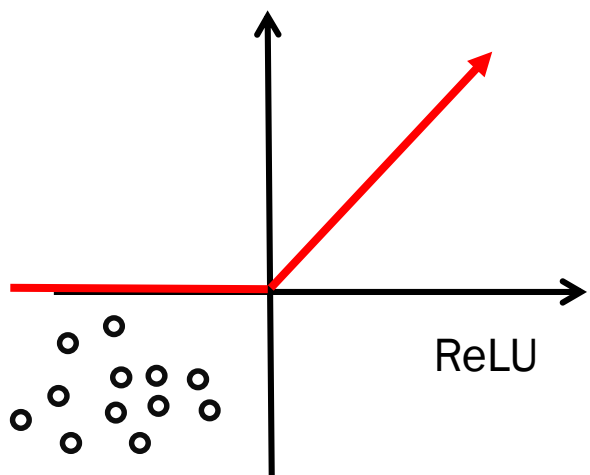
Models are sensitive to weight initialization
Keeping the shape of initialization valid to initiate parameters with better values
randomizing the initial weights, so that the inputs of each activation function fall
within the sweet range of the activation function.
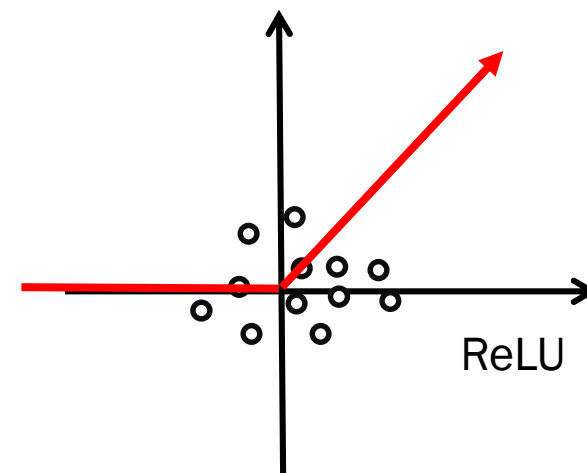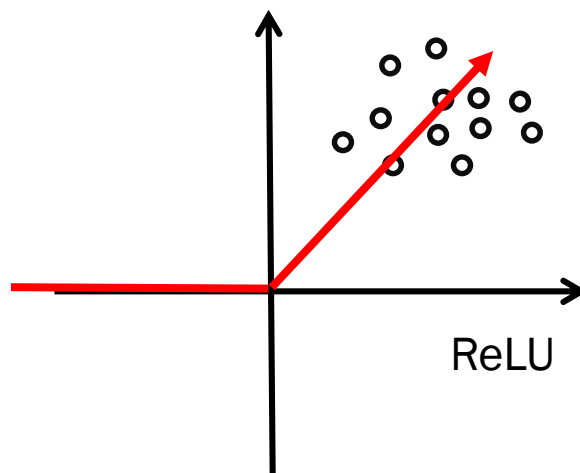Ideally, none of the neurons should start with a trapped situation.

# Batch normalization

Normalize distribution of each input feature in each layer across each minibatch to Normal distribution $N\sim(0,1)$
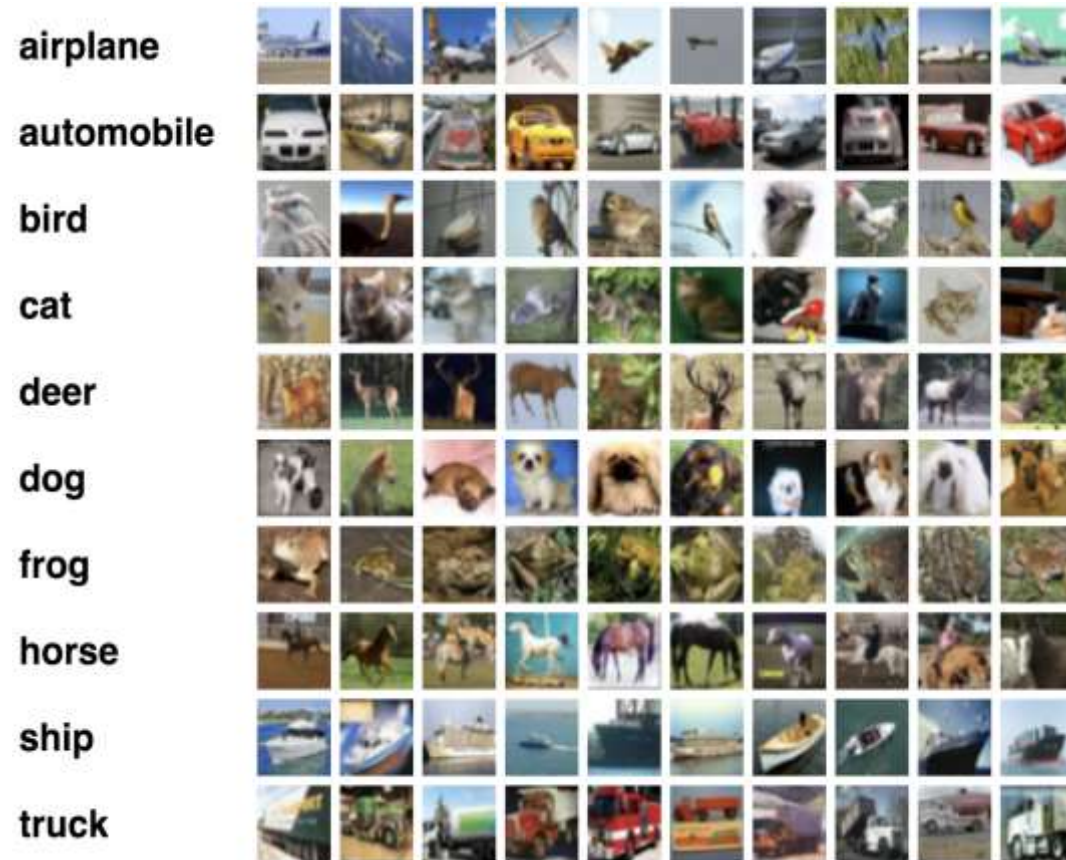
ReLU

ReLU

ReLU

Nothing happens!

Better(greater) learning rate
Faster convergence

# CIFAR10



- 60K 32x32 colour images

- 10 classes (6 K images per class)

- 50K training images

- 10K test images

https://www.cs.toronto.edu/~kriz/cifar.html

# *Lab 4B: CIFAR10 classification - CNN*

| Exercise 1:<br>More conv2D layer | # of conv2D |
|---|---|
| Run 1 | 2 |
| Run 2 | 4 |

| Exercise 2:<br>different optimizer | Optimizer |
|---|---|
| Run 1 | RMSprop |
| Run 2 | SGD |

| Exercise 3:<br>different epoch | # of epoch |
|---|---|
| Run 1 | 10 |
| Run 2 | 30 |

**You may want to try**

1. Use different cost function

2. Use different learning rate

   https://pytorch.org/docs/stable/optim.html?highlight=optimizer#torch.optim.Optimizer

**Break room**

**Session ends:**

# Thank you!