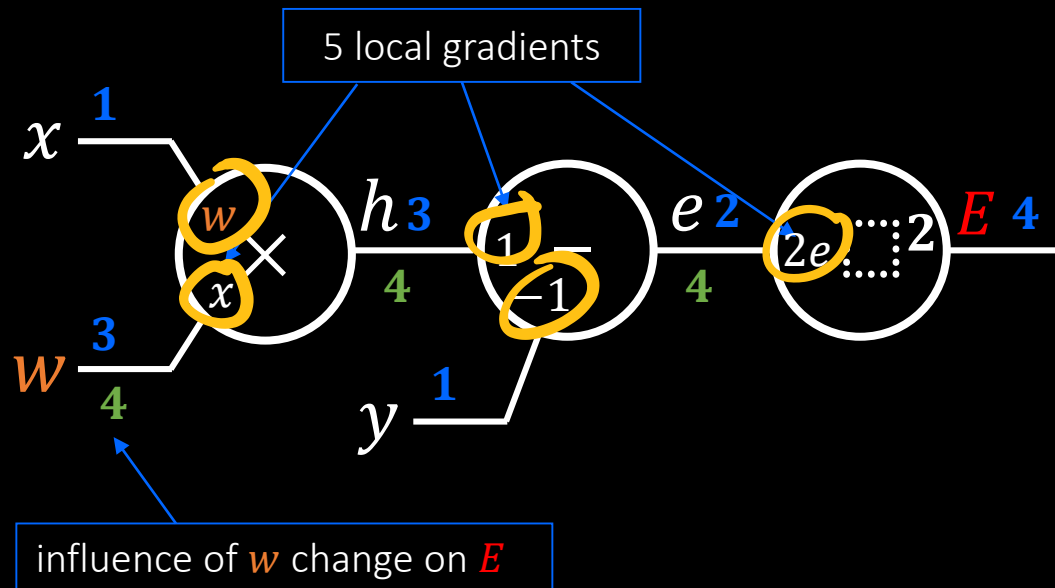AI and Deep Learning

# Deep Learning

Jeju National University

Yungcheol Byun

# Agenda

- Merging gates in a computation graph
- Vanishing gradient and ReLU
- MNIST application
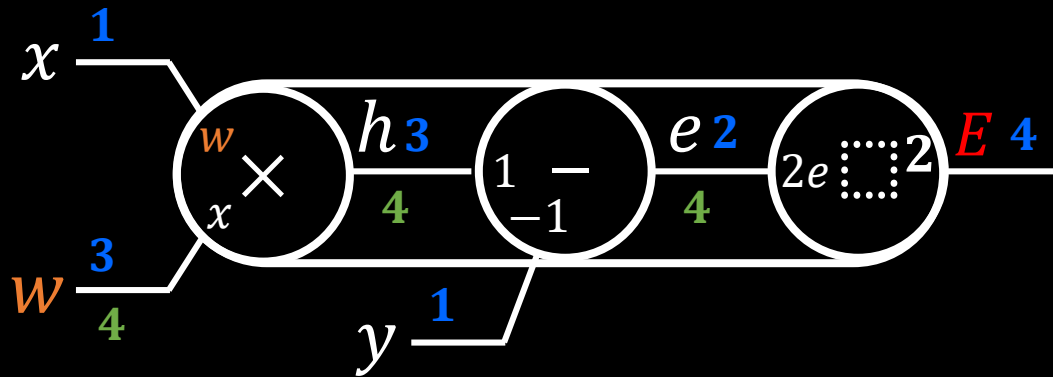- Overfitting and drop-out
- Deep Learning

# Influence of $w$ change on $E$



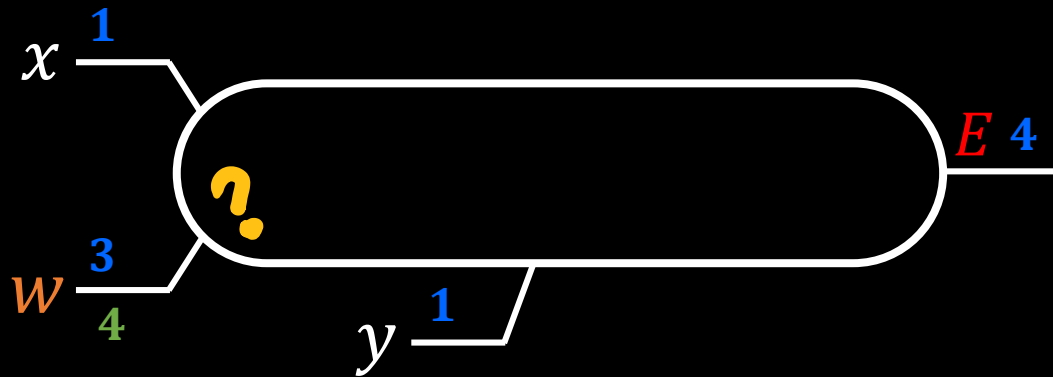is multiplication of all the local gradients in the graph (chain rule)
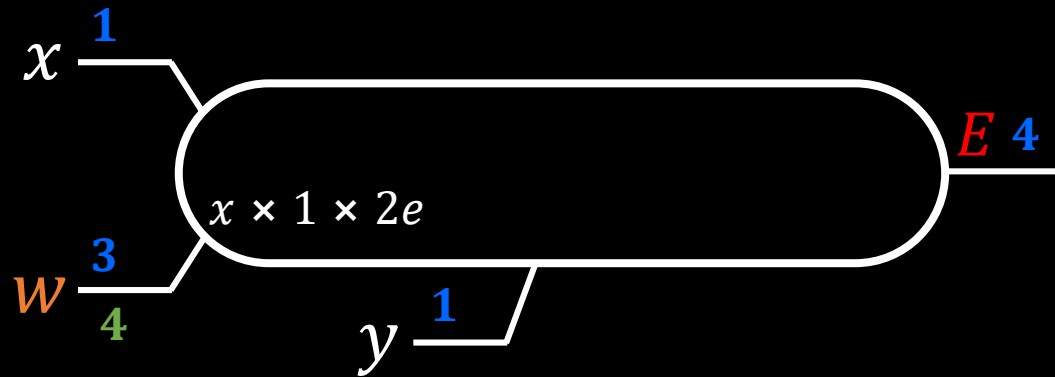
# Influence of $w$ change on $E$

Merging gates

# Influence of $w$ change on $E$

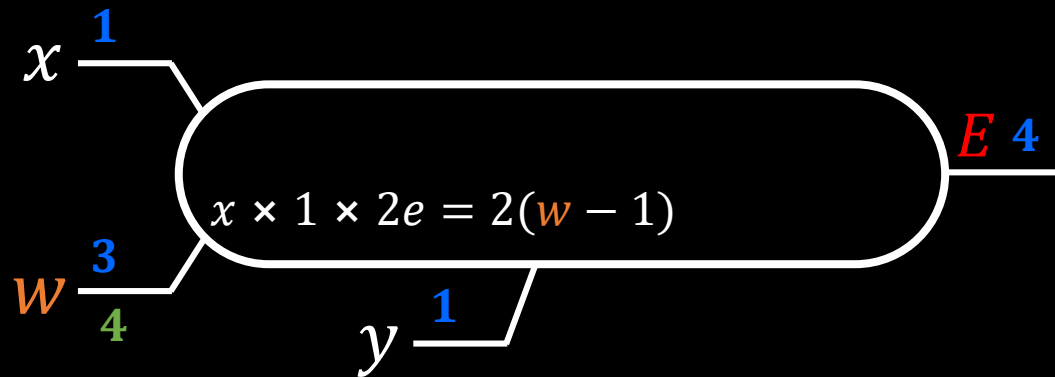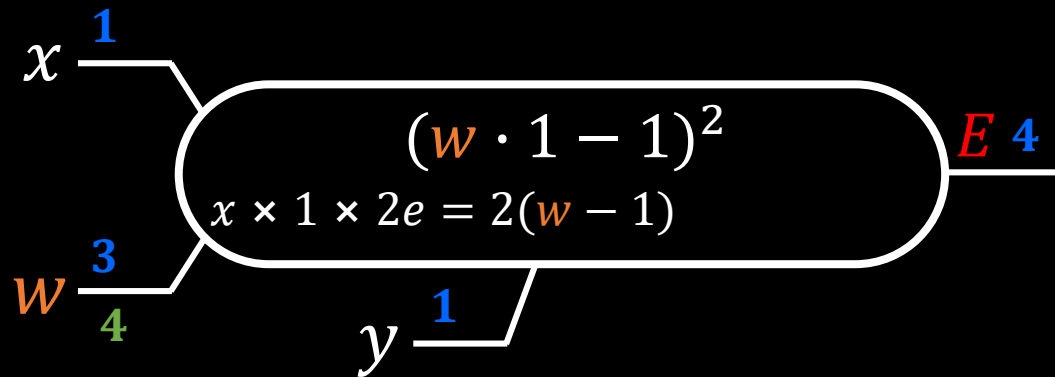Merging gates

# Influence of $w$ change on $E$

Merging gates

$x$ —— 1

$E$ 4

$x \times 1 \times 2e$

$w$ —— 3 / 4

$y$ —— 1

is multiplication of all the local gradients in the graph (chain rule)

# Influence of $w$ change on $E$

Merging gates



$x \times 1 \times 2e = 2(w-1)$

is multiplication of all the local gradients in the graph (chain rule)

# Influence of $w$ change on $E$

Merging gates

$x$ **1**

$(w \cdot 1 - 1)^2$

$x \times 1 \times 2e = 2(w - 1)$

$E$ **4**

$w$ **3** **4**

$y$ **1**

Therefore, the local gradient is derivative of the function.
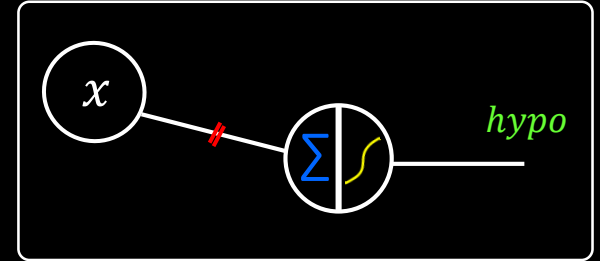
# Influence of $w$ change on $E$

Derivative of $E$ with respect to $w$

$$E = (w \cdot 1 - 1)^2$$

$$\frac{\partial E}{\partial w} = \frac{\partial}{\partial w}(w \cdot 1 - 1)^2 = 2(w - 1)$$
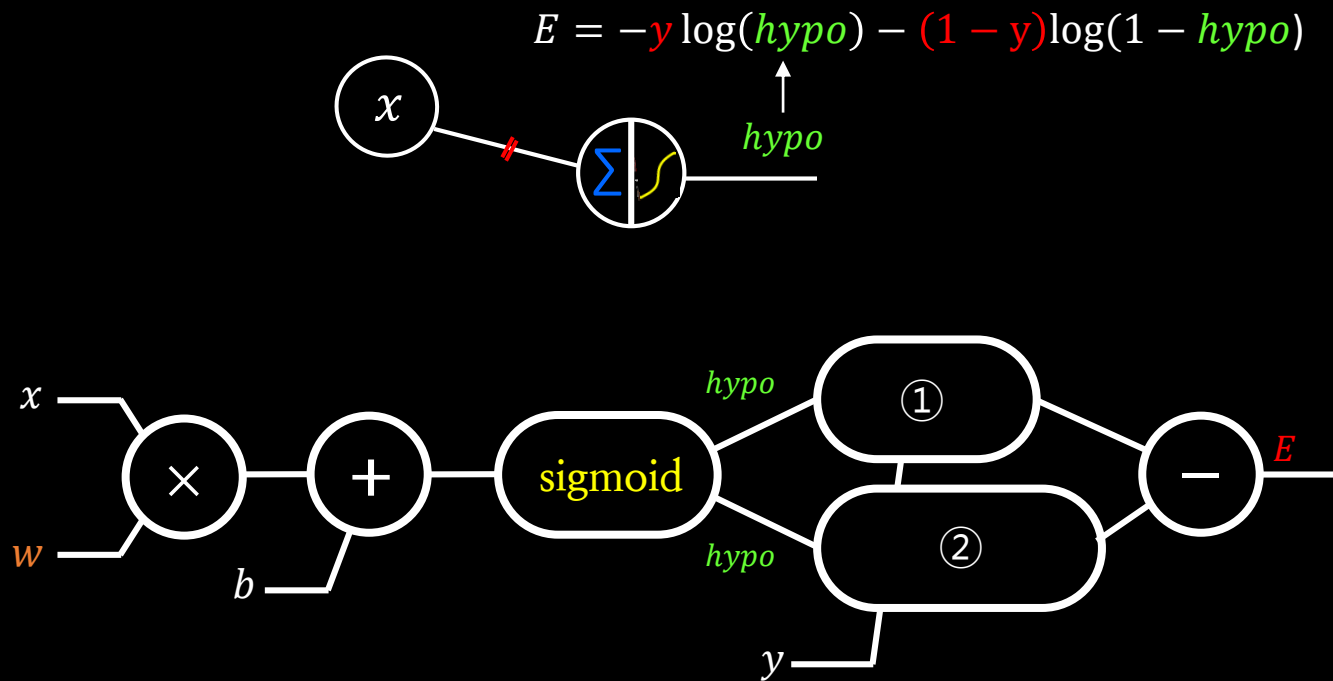
# Cost/Error function

for logistic regression

$$hypo = \frac{1}{1 + e^{-wx}}$$

$$E = -y \log(hypo) - (1 - y)\log(1 - hypo)$$

# Computational Graph

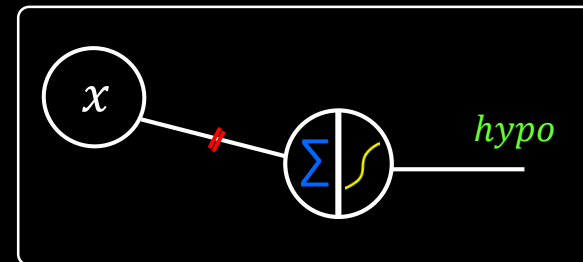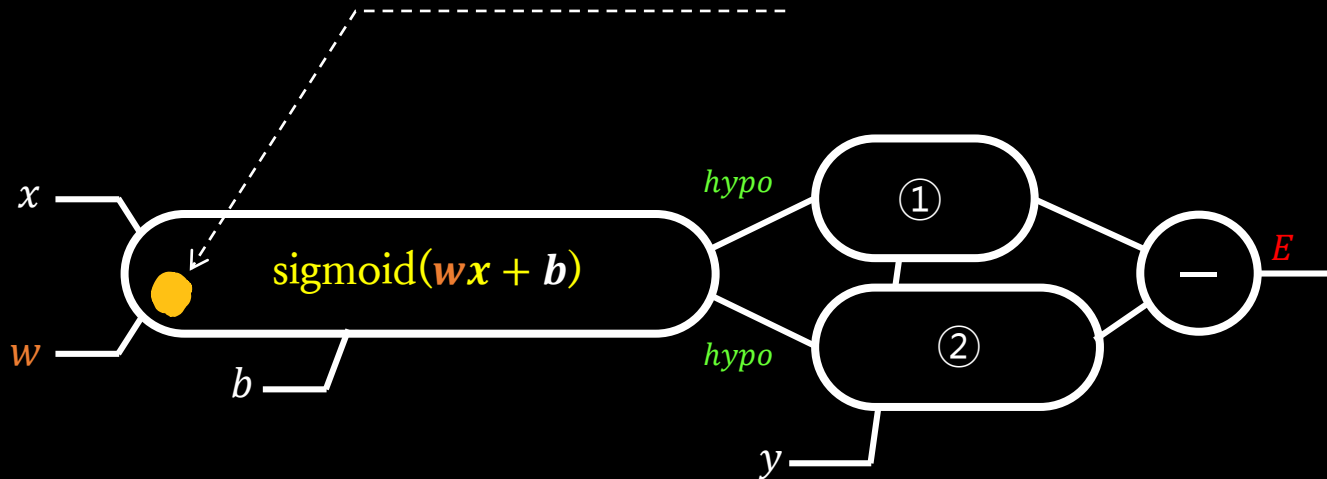$$E = -y \log(hypo) - (1-y)\log(1 - hypo)$$

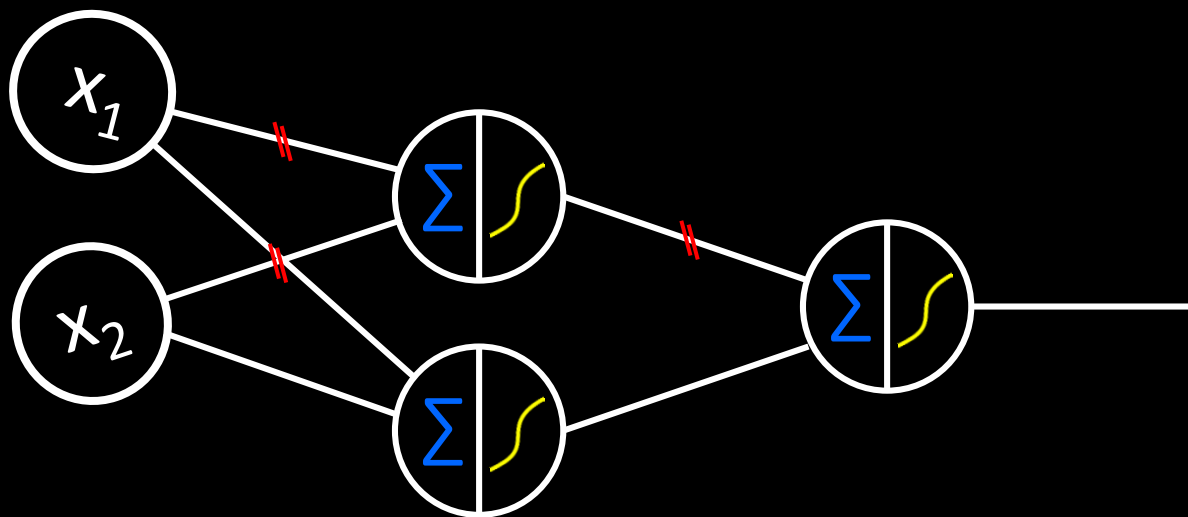

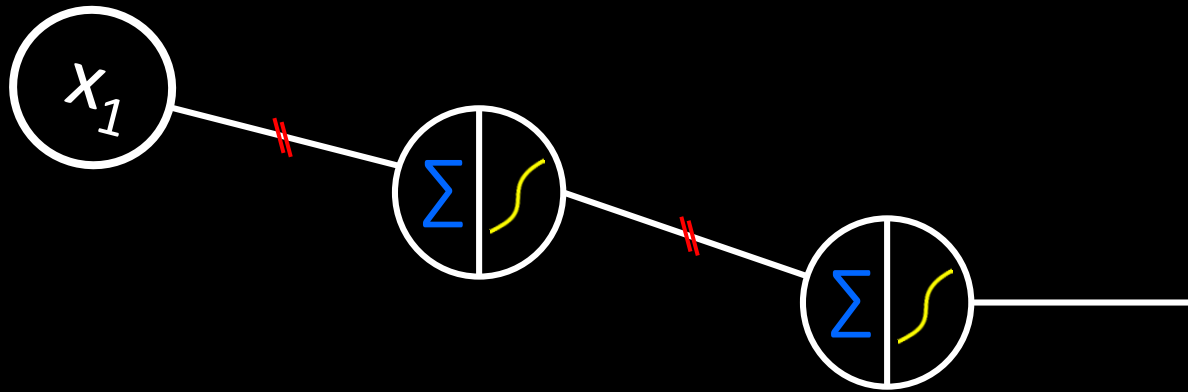$$\frac{\partial E}{\partial w} =$$

# Computational Graph

Merging gates
How we get the local gradient of the merged gate?

# 3-layer NN

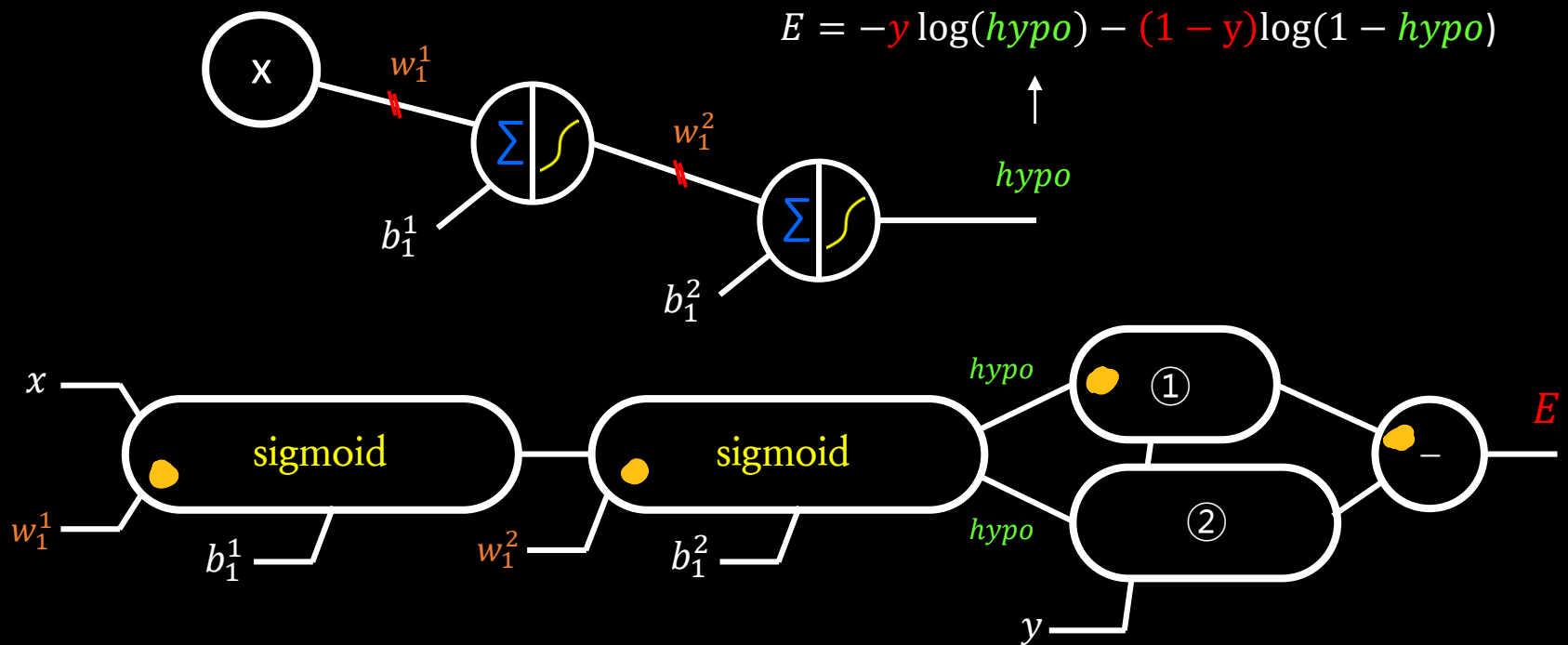# 3-layer NN (simplified)
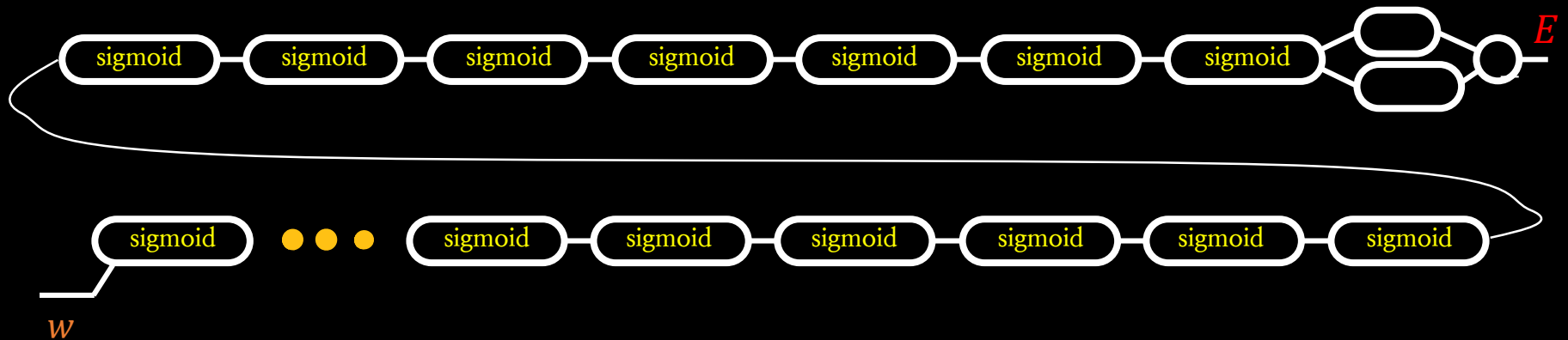
# Influence of $w$ change on $E$



$$E = -y \log(hypo) - (1-y)\log(1 - hypo)$$

# 10-layer Neural Network

The giant moster, computational graph!

# Influence of $w$ change on $E$



$$\frac{\partial E}{\partial w} = \; ?$$

*Hint*: chain rule!

# Vanishing Gradient

- The derivative of sigmoid function is (1-sigmoid) * sigmoid

- Two multiplication of sigmoid for a single neuron, 20 multiplications for 10 connected neurons

- Each sigmoid squashes the input value into the value between 0 and 1.

# Vanishing Gradient

- The influence of $w$ change on $E$ is *many* multiplications of the values between 0 and 1, which gives us almost 0.
- Vanishing Gradient
- $w = w - \alpha \cdot$ (almost 0)
- $b = b - \alpha \cdot$ (almost 0)
- Therefore, no change in $w$ and $b$

# (Lab) 18.py

- XOR problem using 4-layer neural networks
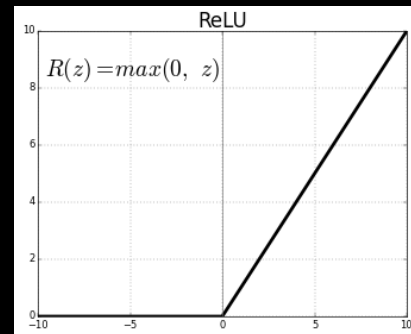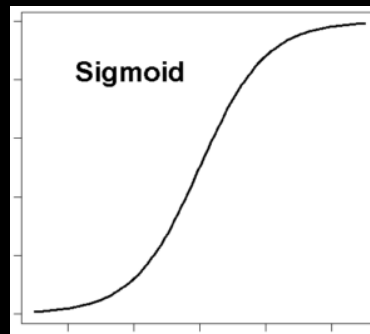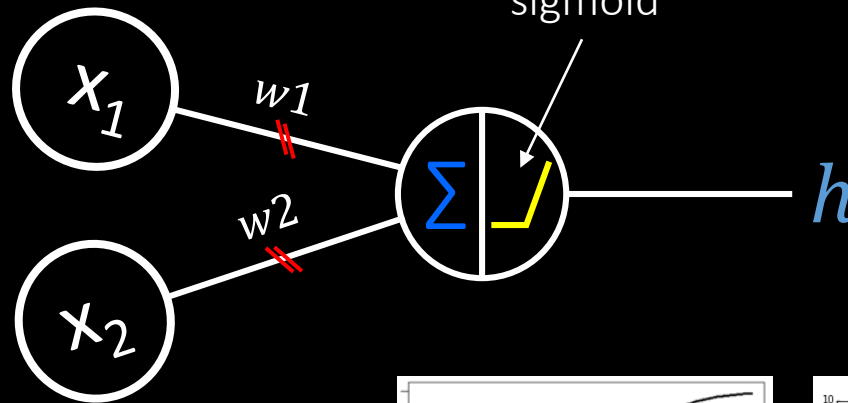- Failed owing to vanishing gradient

# The Dark Age in Artificial Intelligence and Neural Networks (~2006)

since back-propagation by Hinton in 1986

# ReLU

using ReLU (Rectified Linear Unit) activation function instead of sigmoid

$x_1$

$w1$

$x_2$

$w2$

$\Sigma$

$h$

Sigmoid

ReLU

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

# (Lab) 19.py

- Solving vanishing gradient problem using ReLU activation function

- Back-propagation is working by using ReLU.
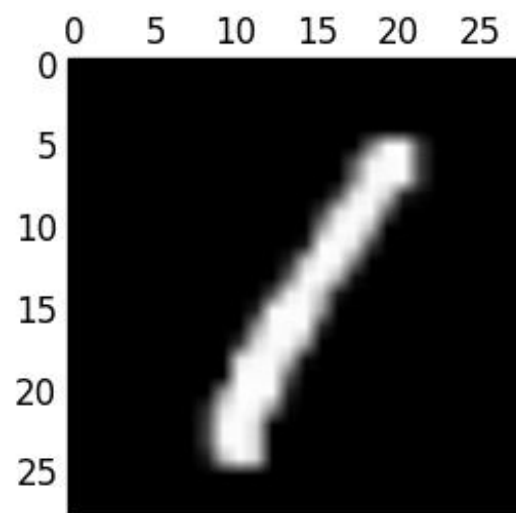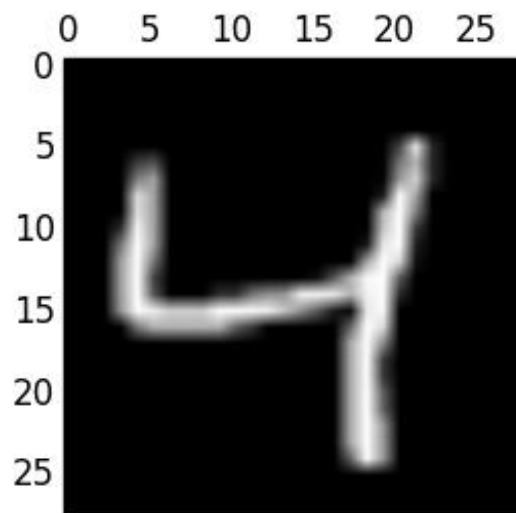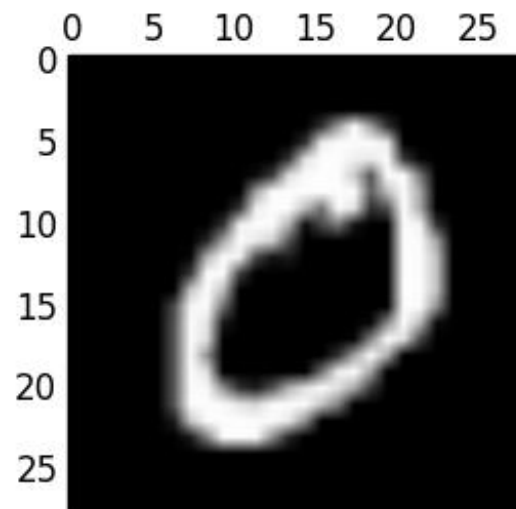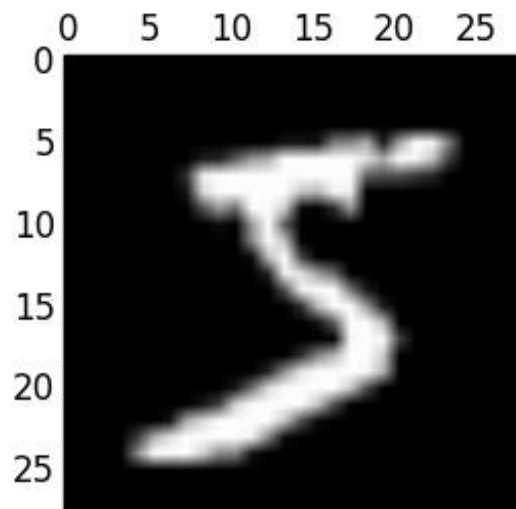
So, now can go deeper.

# MNIST



Modified National Institute of
Standards and Technology
(USA)

# MNIST

# (Lab) 20.py

- 60,000 training images + 10,000 testing images
- Input image : 28 * 28 pixels → 784 pixels
- 784 dimension
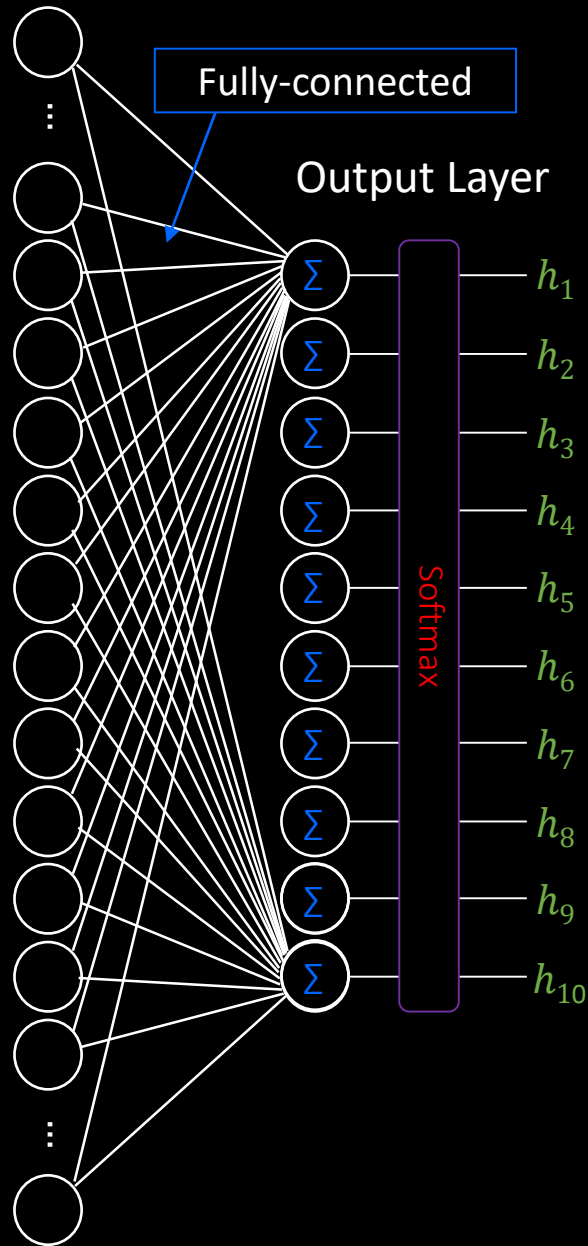- 10 classes (output: 0 ~ 9)
- Softmax
- 90.23% of recognition rate

Input Layer

Fully-connected
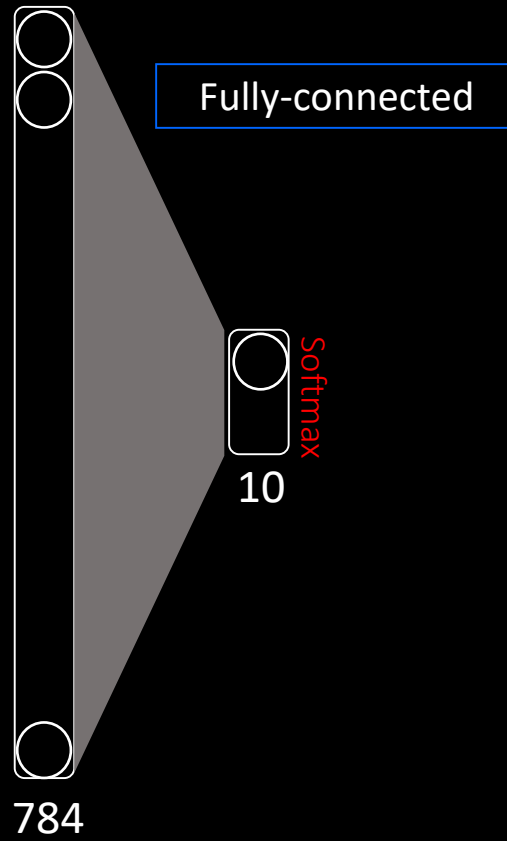
Output Layer

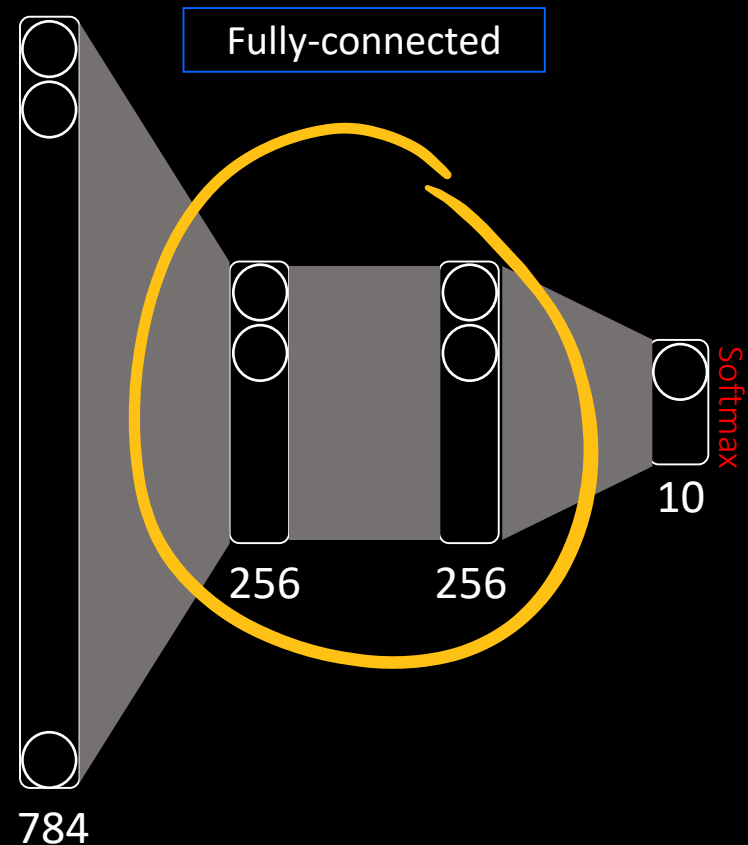$\Sigma$ — $h_1$
$\Sigma$ — $h_2$
$\Sigma$ — $h_3$
$\Sigma$ — $h_4$
Softmax
$\Sigma$ — $h_5$
$\Sigma$ — $h_6$
$\Sigma$ — $h_7$
$\Sigma$ — $h_8$
$\Sigma$ — $h_9$
$\Sigma$ — $h_{10}$

784

10

Fully-connected

Softmax

10

784

# (Lab) 21.py

- Deep Neural Network (4-layer)
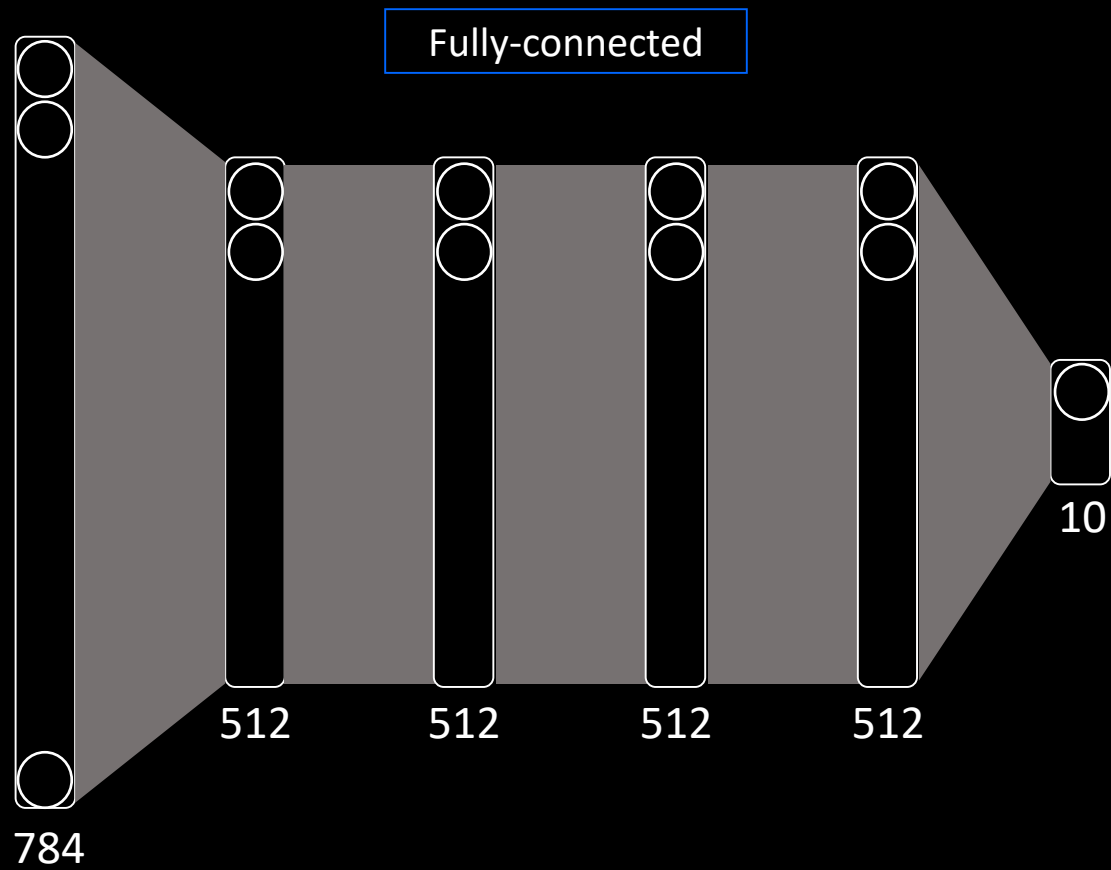- ReLU
- 94.55% accuracy

Fully-connected

784

256

256

10

Softmax

# (Lab) 22.py

- Applying initialization method for $w$ and $b$, not randomly
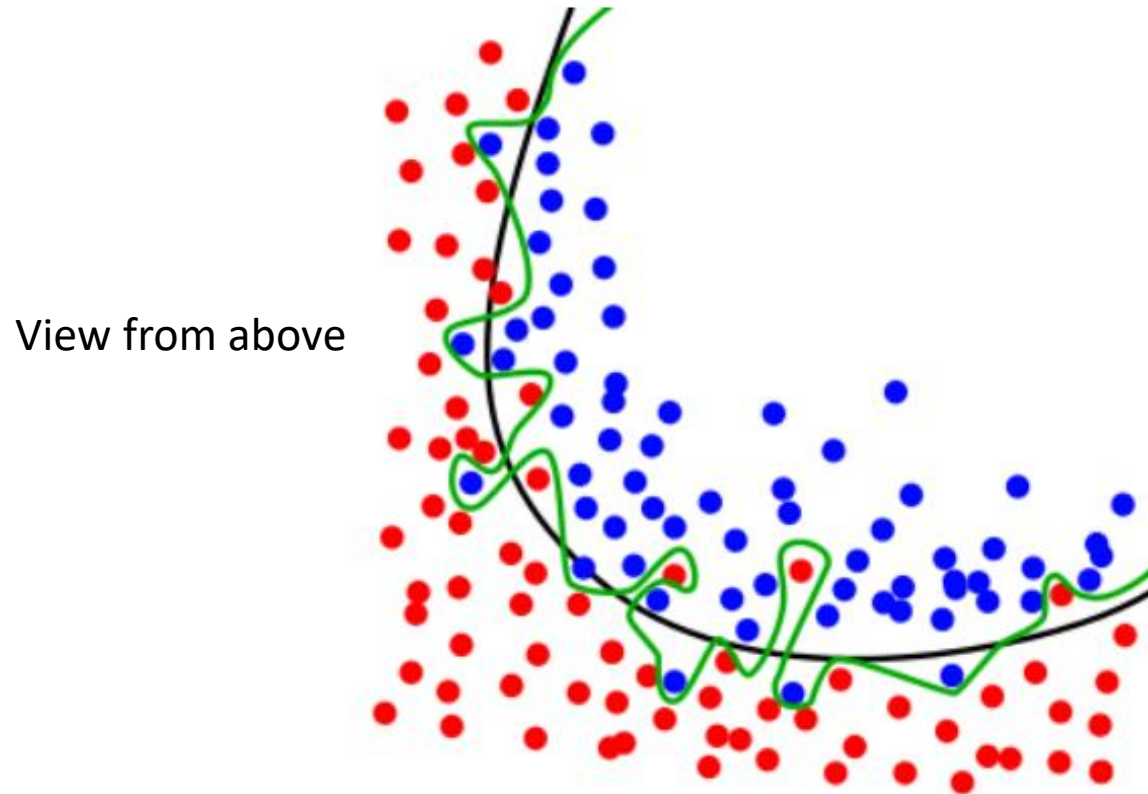
- 97.23% of accuracy

# (Lab) 23.py

- Applying initialization method for $w$ and $b$, not randomly

- 6-layer deep neural networks

- 97.83% of accuracy

# Overfitting and drop-out

- The deeper the network is, the more the decision boundary is complex.

- Good at learning data but errors for testing data → overfitted to learning data

- Making it less complex by drop-out some neurons while learning.
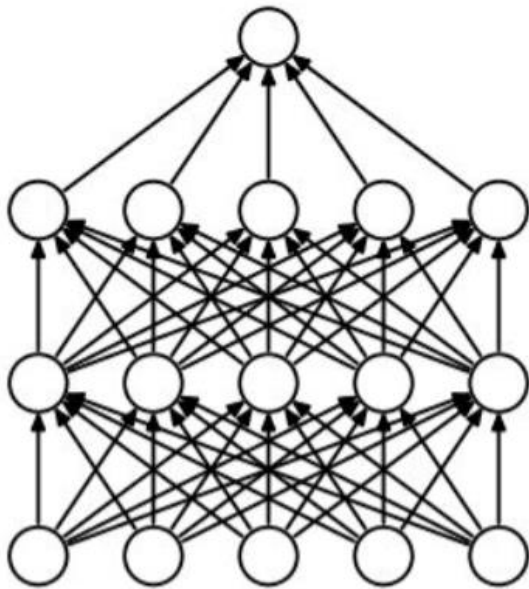
# Which do you think is desirable decision boundary?

View from above



*While the black line fits the data well, the green line is overfit.*
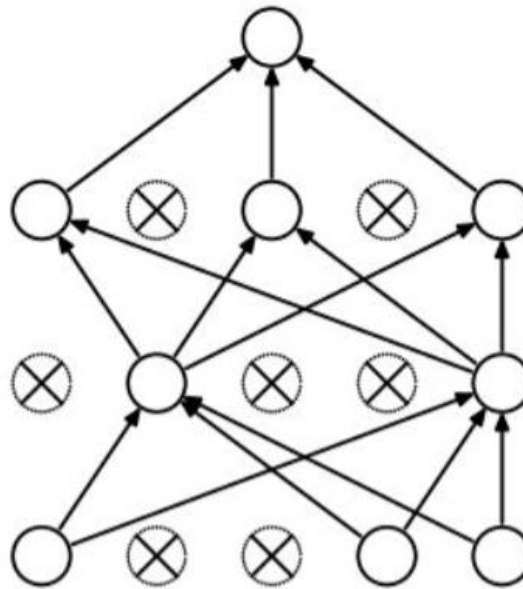
https://elitedatascience.com

# Regularization: **Dropout**
"randomly set some neurons to zero in the forward pass"



(a) Standard Neural Net     (b) After applying dropout.
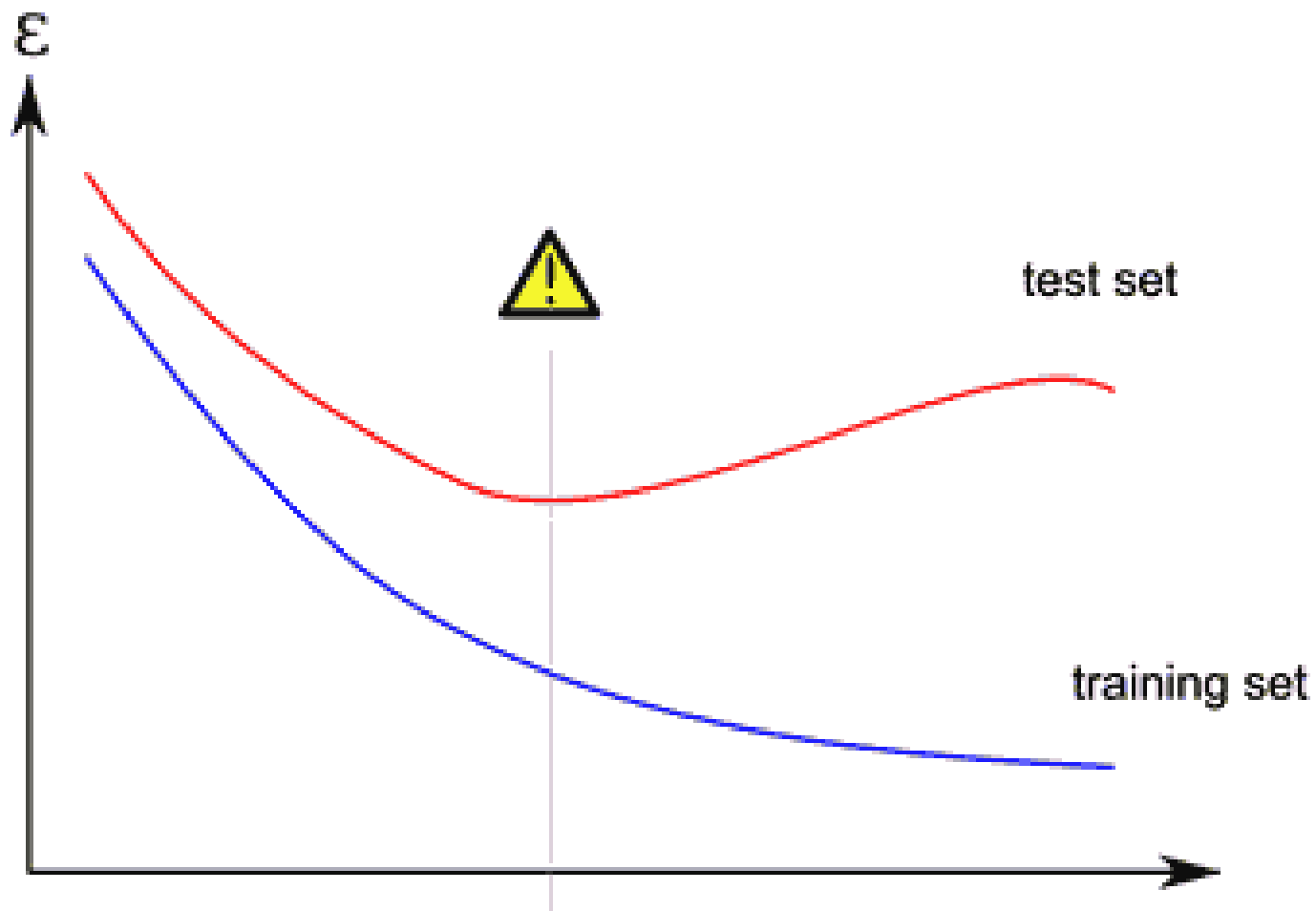
*[Srivastava et al., 2014]*

# (Lab) 24.py

- Applying dropout
- 98.13% of recognition accuracy
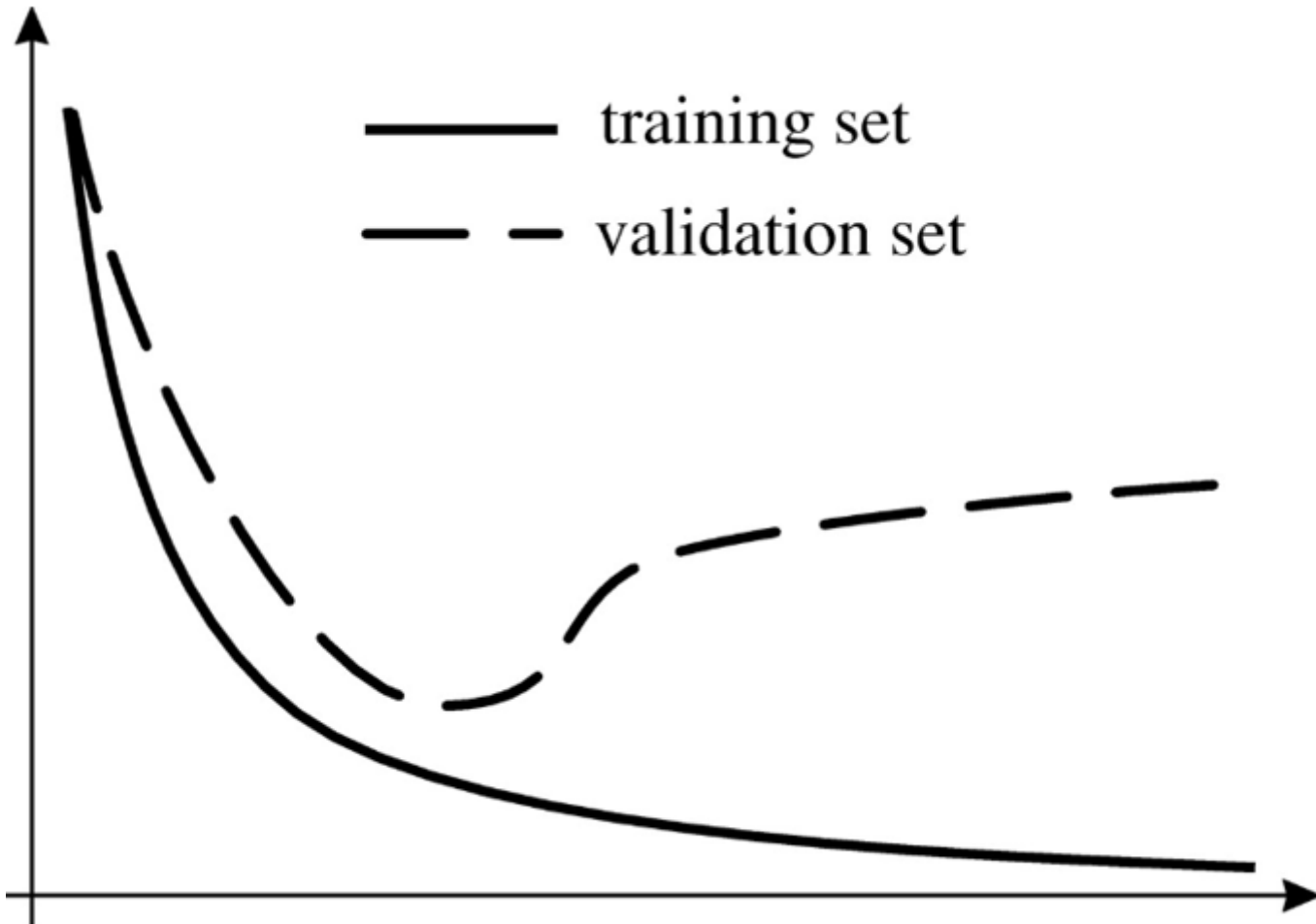
# How to Prevent Overfitting

- Train with more data
- Reduce features
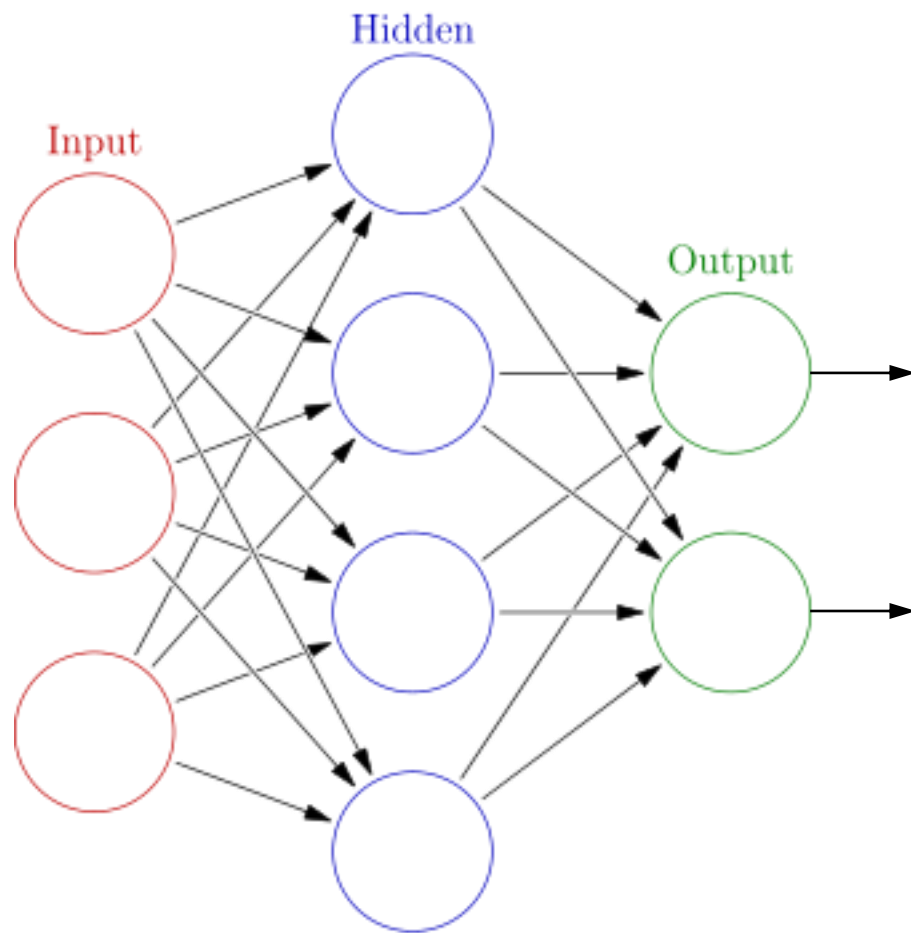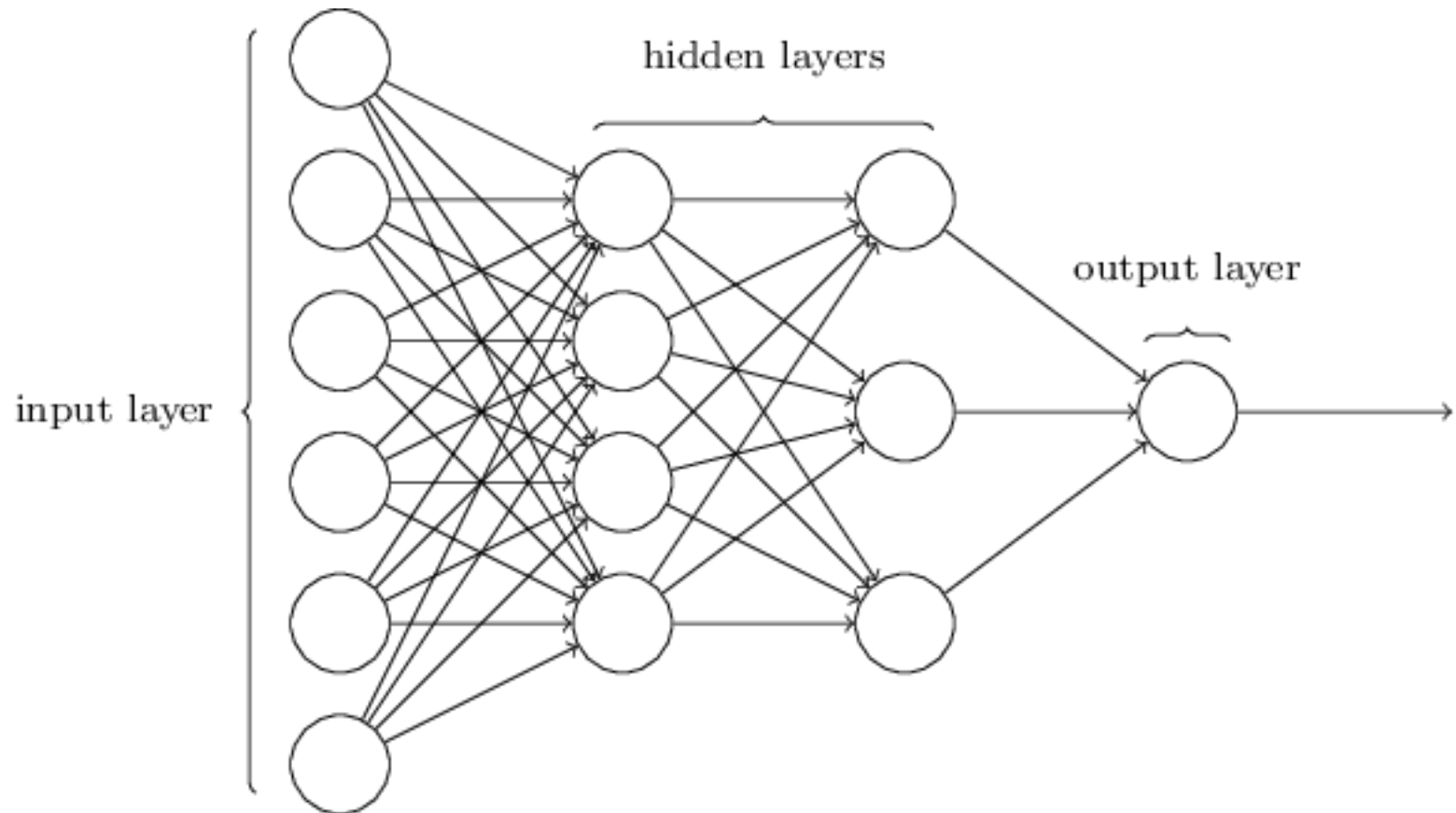- Early stopping
- Ensemble
- regularization

# Early stopping

# Early stopping



Legend:
— training set
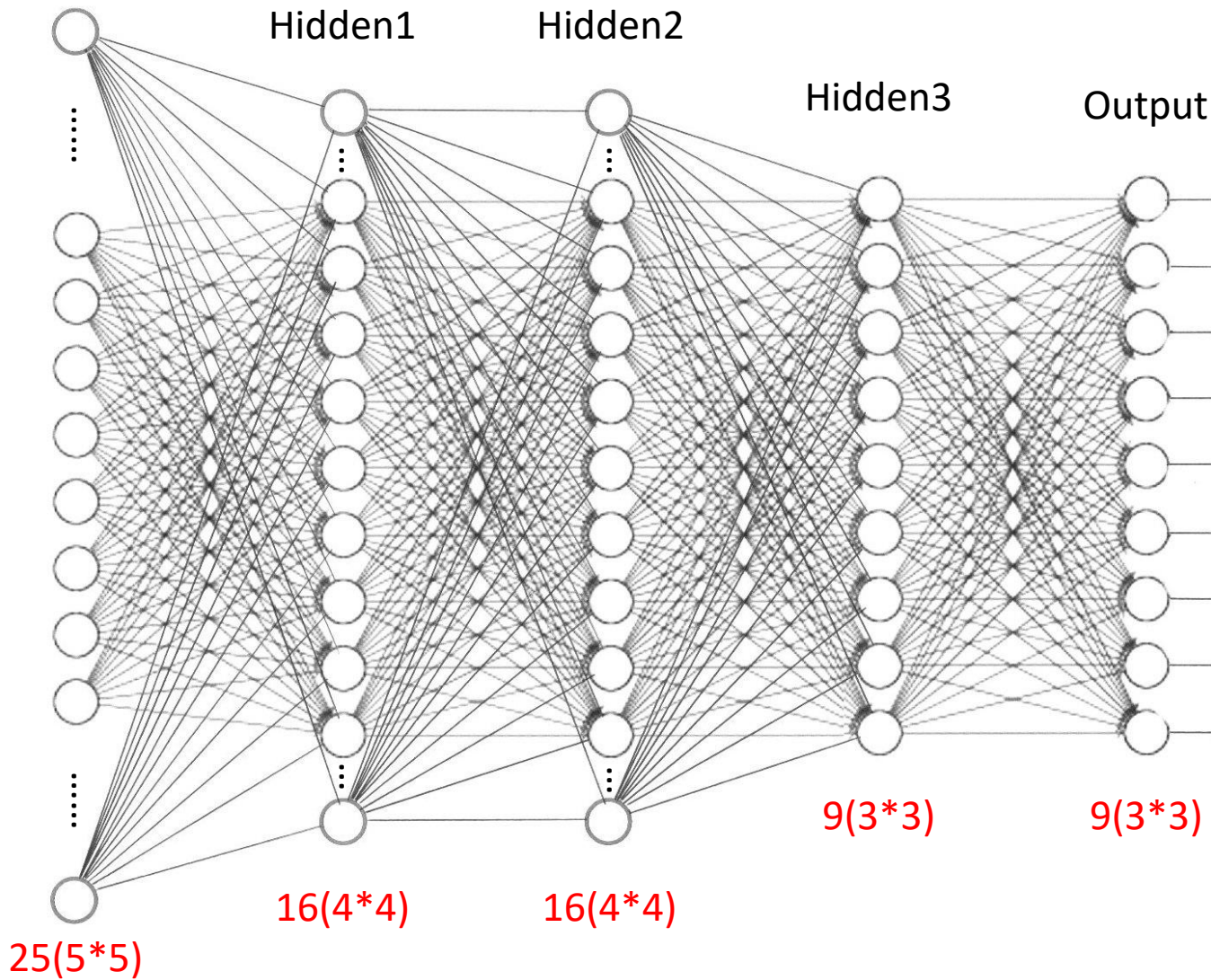– – – validation set

Fully-connected

Hidden

Input

Output

Fully-connected

hidden layers

output layer

input layer

Input

Fully-connected
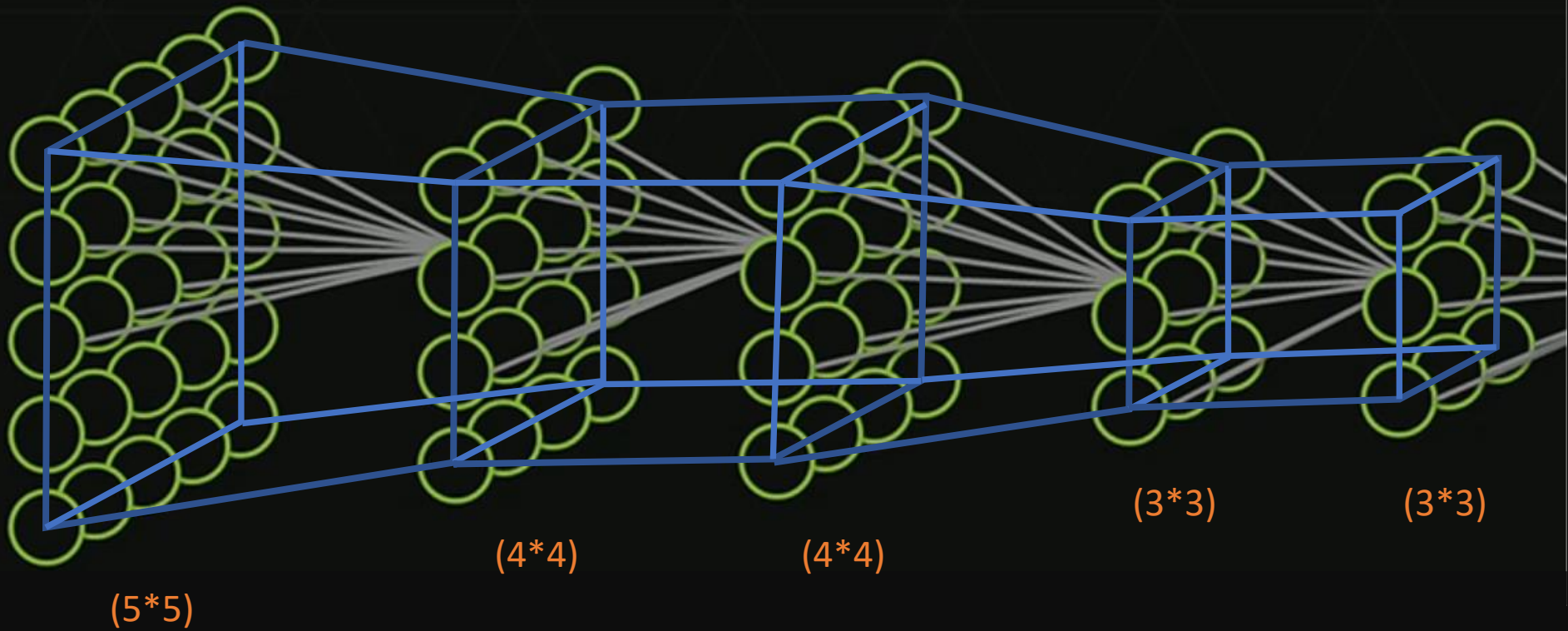
Hidden1　　Hidden2

Hidden3　　Output

9(3*3)　　9(3*3)

16(4*4)　　16(4*4)

25(5*5)

Fully connected, then how many synapses(parameters) are there?
25 * 16 + 16 * 16 + 16 * 9 + 9 * 9 = 881

Fully-connected

(5*5)

(4*4)

(4*4)

(3*3)

(3*3)

Fully connected, so how many connections are there?
25 * 16 + 16 * 16 + 16 * 9 + 9 * 9 = 881

(3*3)   (3*3)

(4*4)   (4*4)

(5*5)

Geoffrey Hinton, Yann LeCun, Yoshua Bengio, Andrew Ng

# Deep Learning

- in early 2000s (2006, 2010, 2012)
- Deep Neural Networks
- Activation functions (ReLU)
- Weight initialization methods
- Dropout (2014)
- Big data
- GPU

Fully-connected

# FCNN

Any problem?