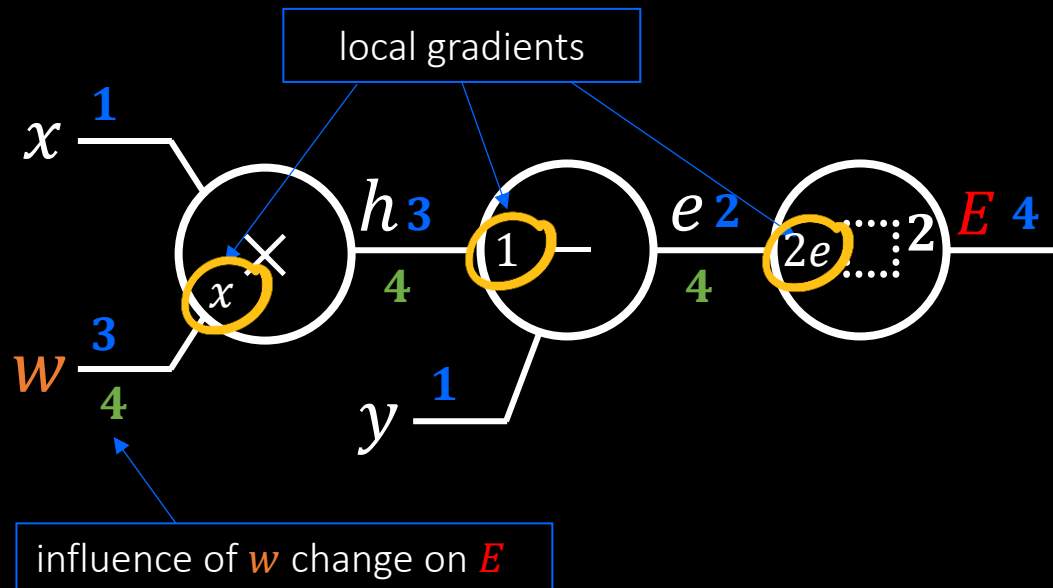AI and Deep Learning

# Deep Learning

Jeju National University

Yungcheol Byun

# Agenda

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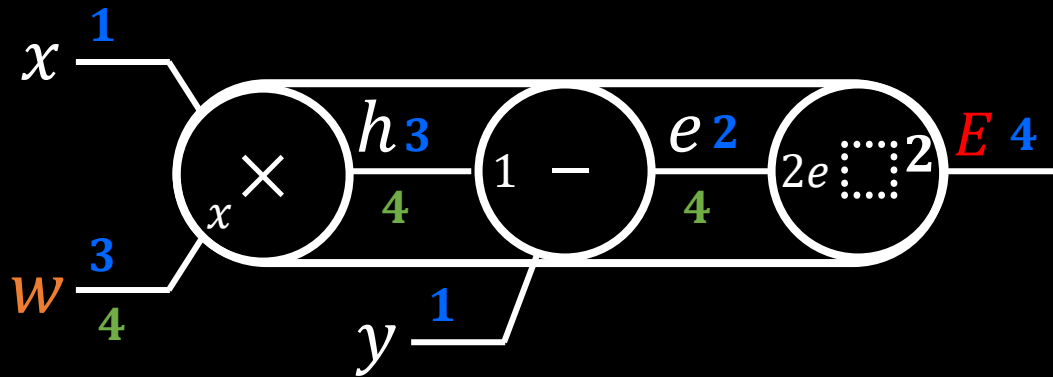- Overfitting and drop-out
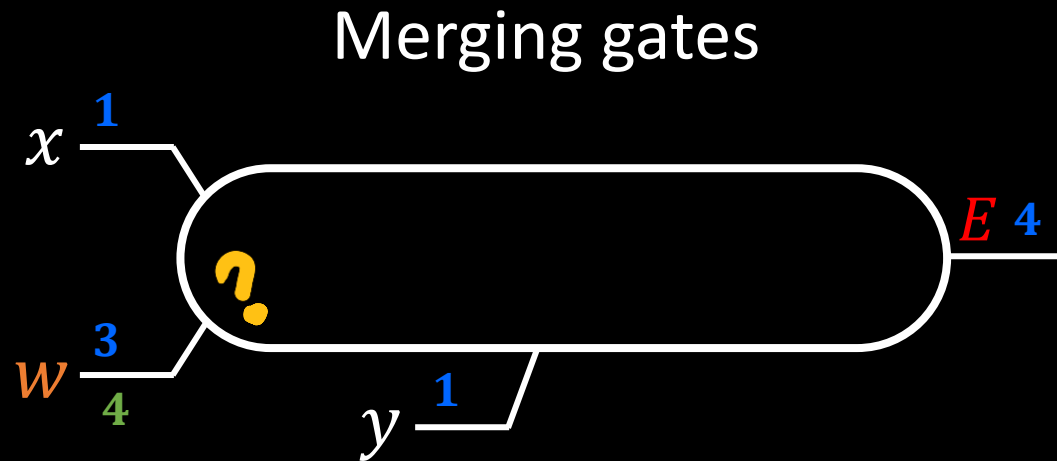- Deep Learning

# Influence of $w$ change on $E$



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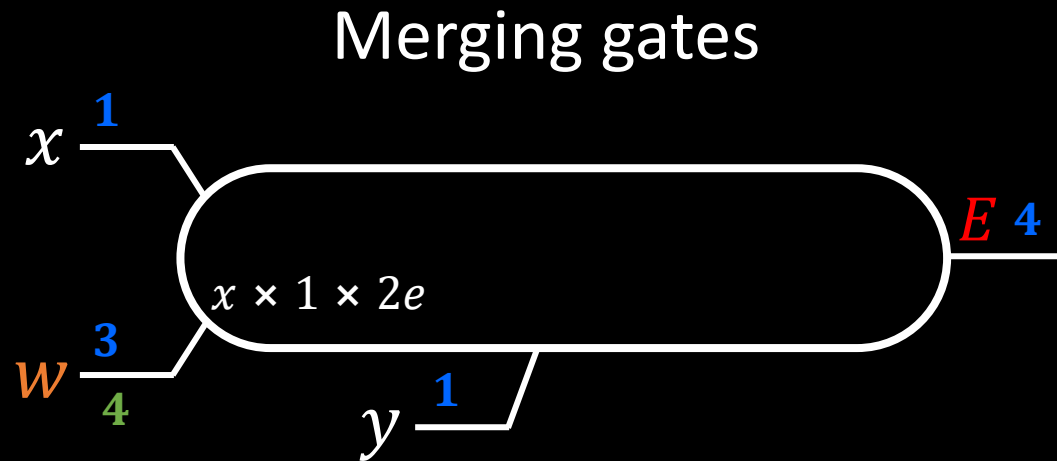
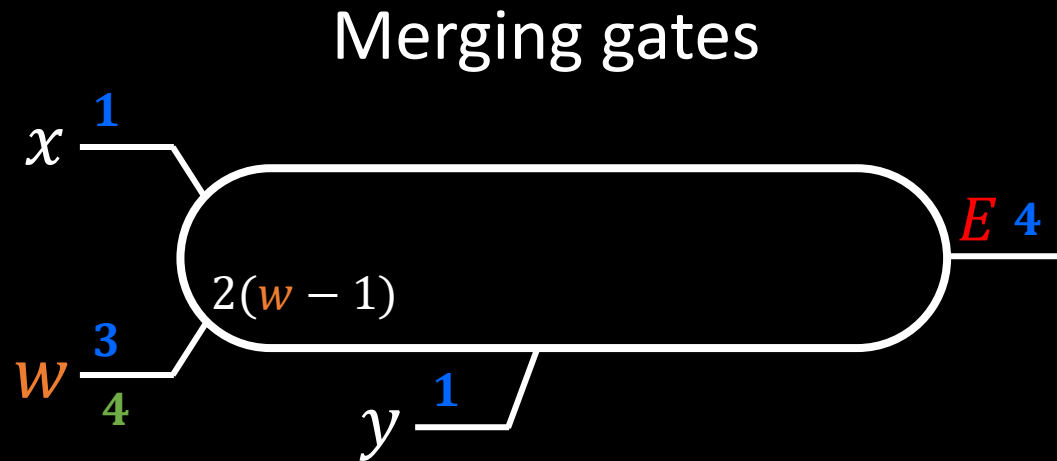# Influence of $w$ change on $E$

# Influence of $w$ change on $E$

Merging gates

# Influence of $w$ change on $E$

Merging gates



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

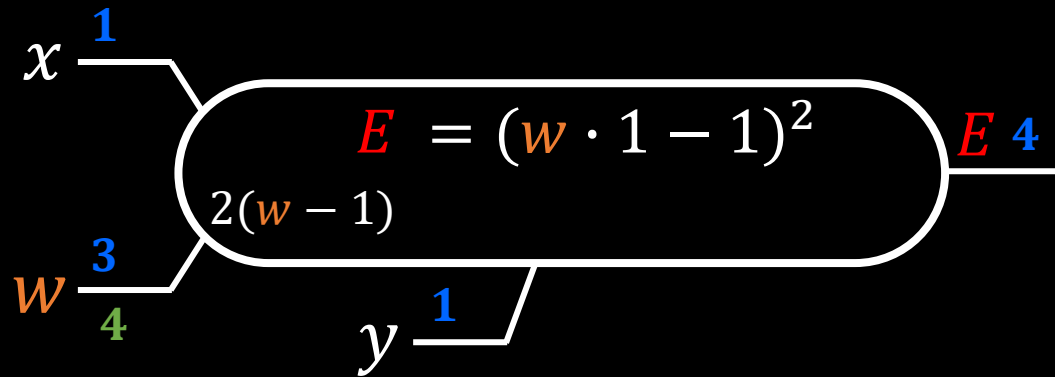# Influence of $w$ change on $E$

Merging gates



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

# Influence of $w$ change on $E$

Merging gates

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

$2(w - 1)$

$x$ — **1**

$w$ — **3** / **4**

$y$ — **1**

$E$ **4**

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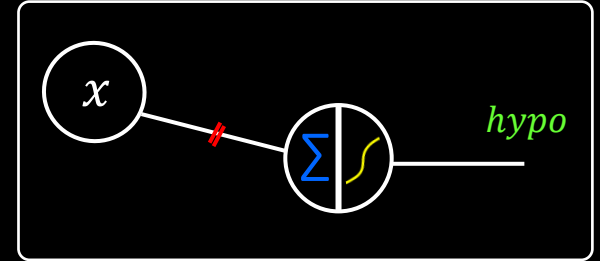
# Influence of $w$ change on $E$

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

Derivative of $E$ with respect to $w$

$$\frac{\partial E}{\partial w} = \frac{\partial}{\partial w}(w \cdot 1 - 1)^2 = 2(w - 1)$$

# Cost/Error function
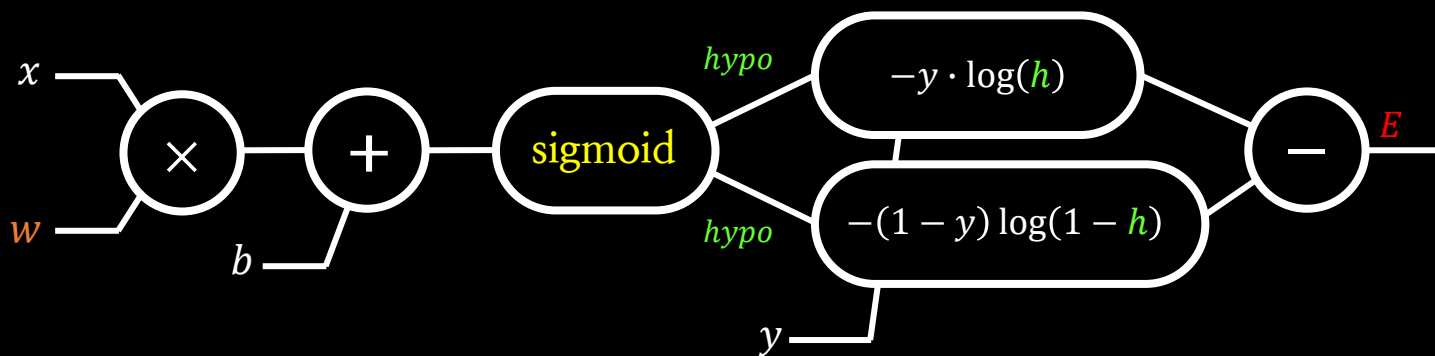
for logistic regression
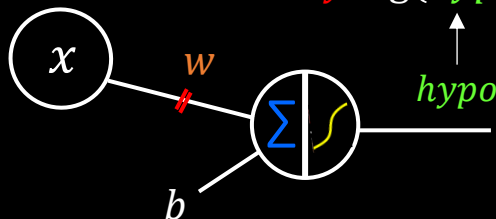
$$hypo = \frac{1}{1 + e^{-wx}}$$

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

Binary Cross Entropy

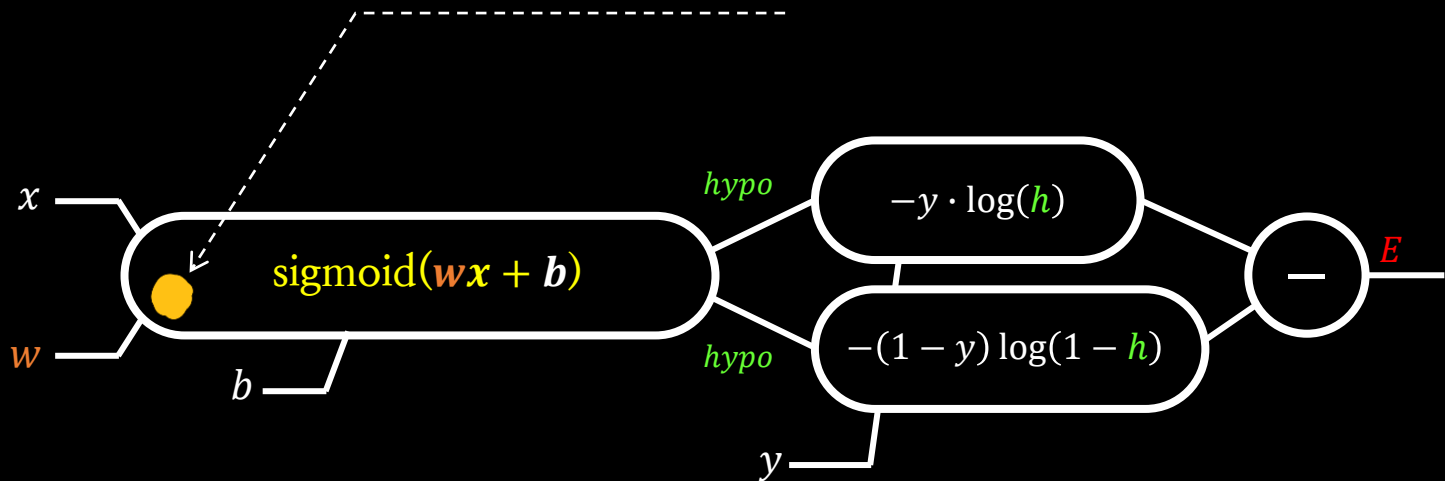# Computational Graph

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



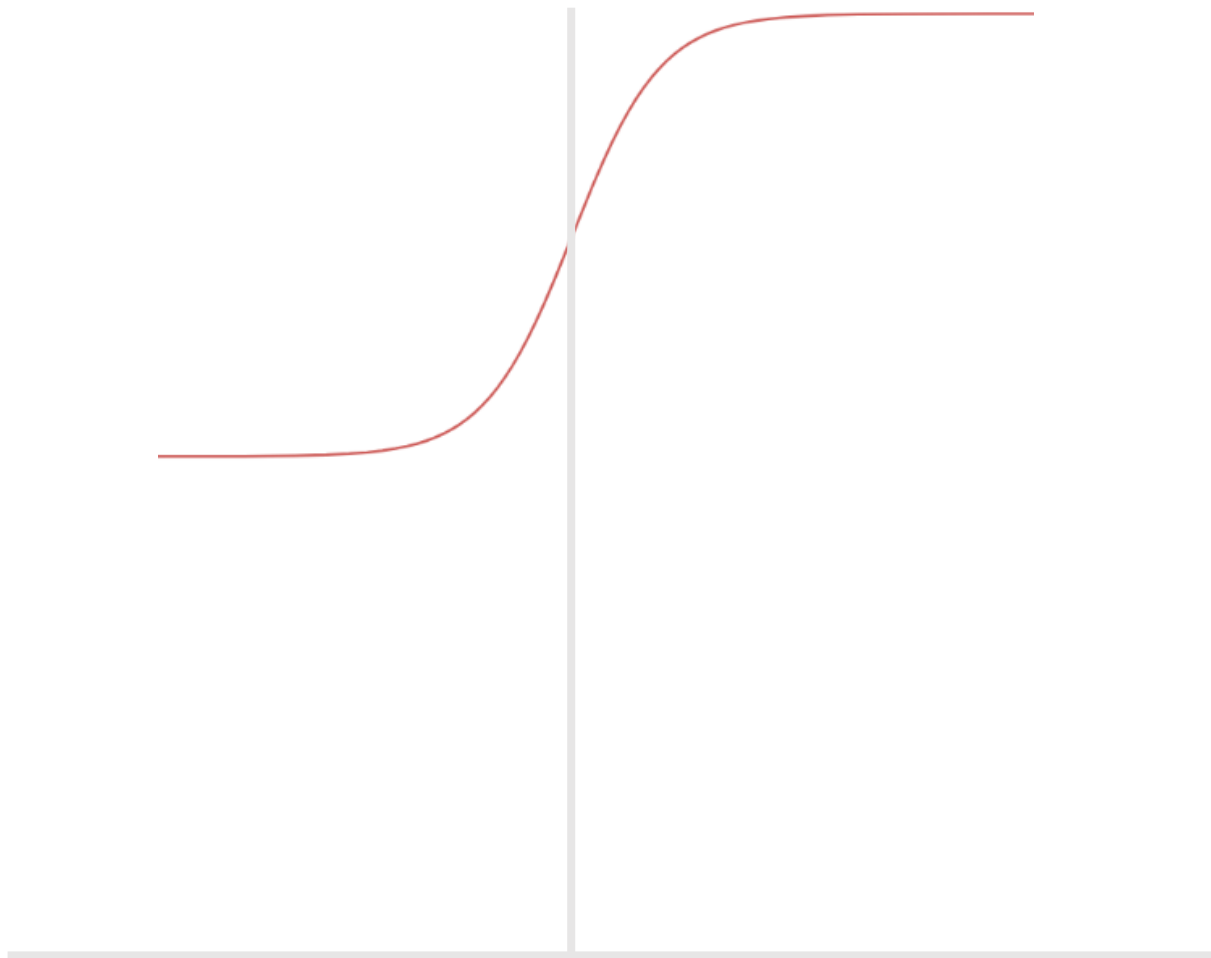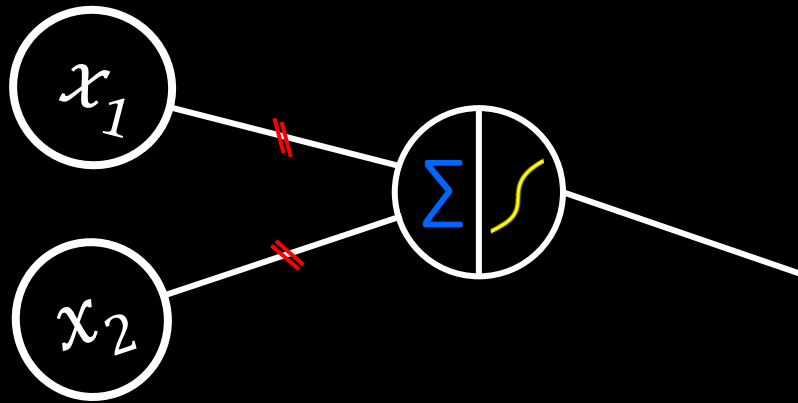$$\frac{\partial E}{\partial w} =$$

# Computational Graph

Merging gates

How we get the local gradient of the merged gate(sigmoid)?

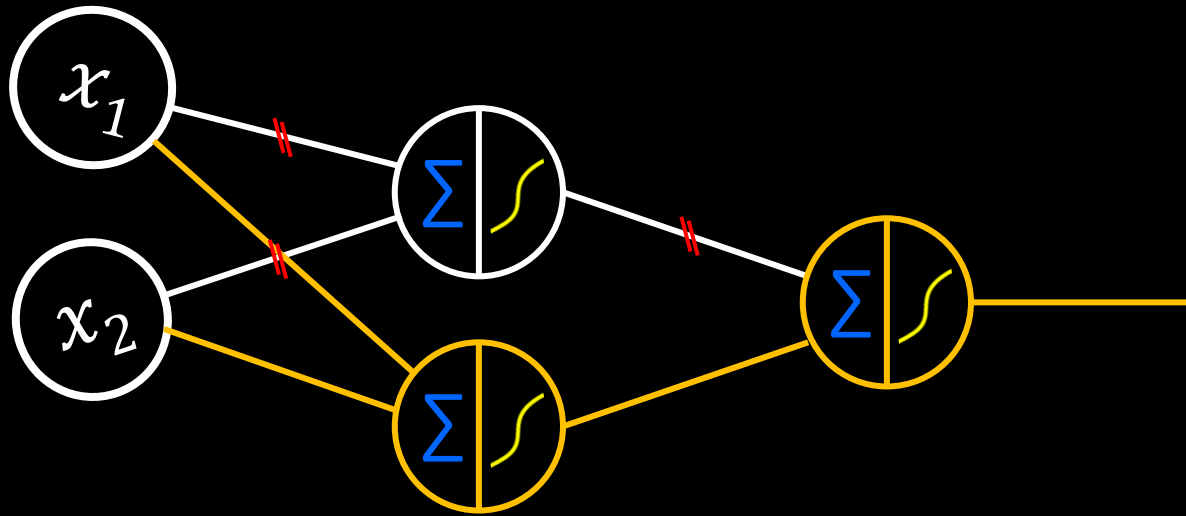$$\frac{\partial sigmoid}{\partial w} = sigmoid(1 - sigmoid)$$

# 1 Neuron(2-layer)

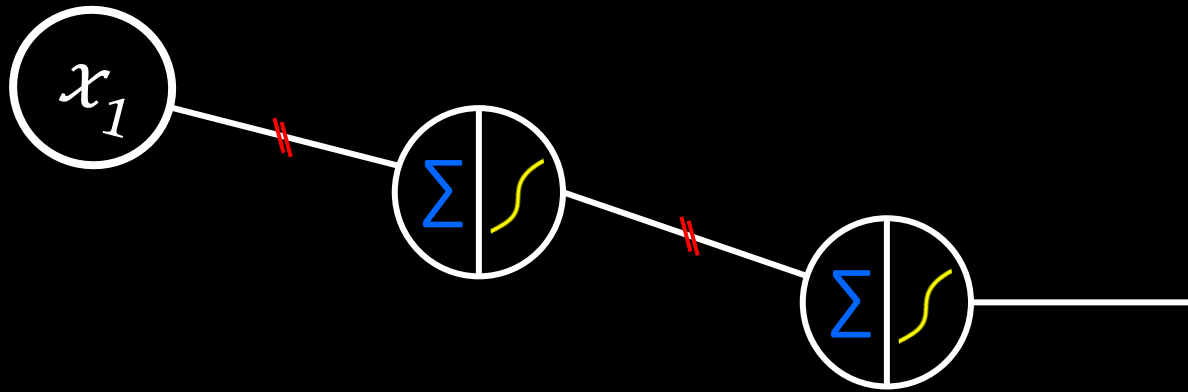# 3-layer NN



Why? what for?

# 3-layer NN (simplified)

# Influence of $w$ change on $E$

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

①  ②

$x$  $w_1^1$  $\Sigma$  $w_1^2$  $\Sigma$  $hypo$

$b_1^1$  $b_1^2$

3-layer NN → 2 sigmoid

$x$  sigmoid  sigmoid  $hypo$ ①  $E$

$w_1^1$  $b_1^1$  $w_1^2$  $b_1^2$  $hypo$ ②

$y$

# 10-layer Neural Network

The giant monster, computational graph!

# Vanishing Gradient

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

- 2 multiplication of sigmoid for a single neuron, 18 multiplications for 9 connected neurons

- Each sigmoid gives us the value between 0 and 1.

Ex) 0.5x0.5 x 0.1x0.9 x 0.8x0.2 x 0.5x0.5 x 0.3x0.7 x 0.4x0.6 x 0.3x0.7 x 0.2x0.8 x 0.9x0.1
= 0.00000010886 x ①

# Vanishing Gradient

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

# (Lab) 18.py

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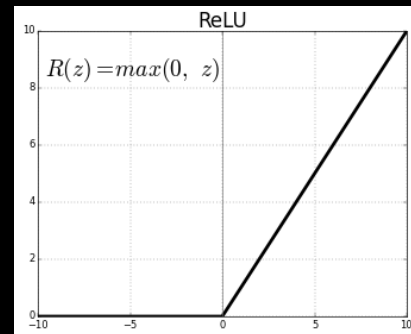- Failed owing to vanishing gradient

https://github.com/yungbyun/neuralnetworks

# The **Dark Age** in AI (~2006)

since **back-propagation** by Hinton in 1986

# ReLU

using ReLU (Rectified Linear Unit) instead of sigmoid

Rectified: 바로잡은, 직각의

$x_1$

$w1$

$w2$

$x_2$

$\Sigma$

$h$

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.

https://github.com/yungbyun/neuralnetworks

So, now can go deeper.

# MNIST
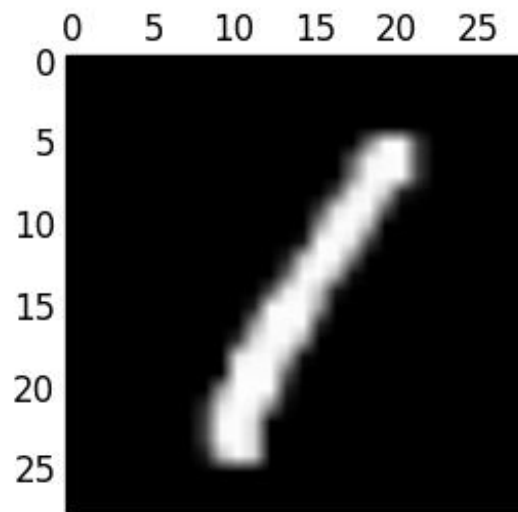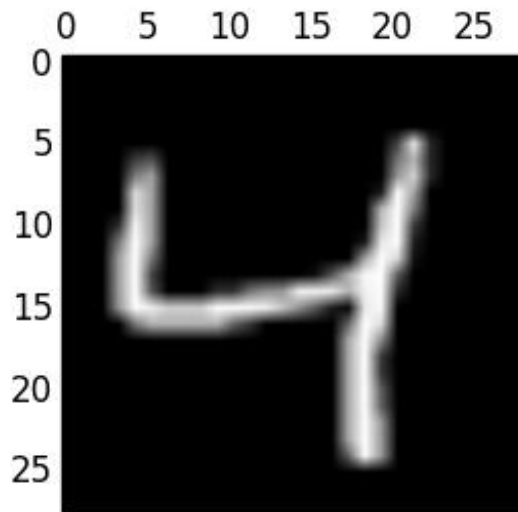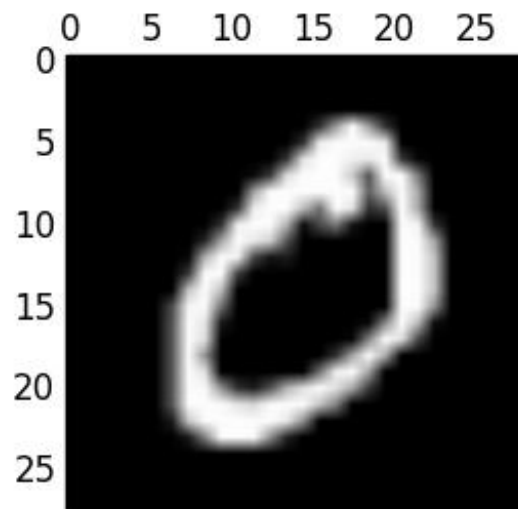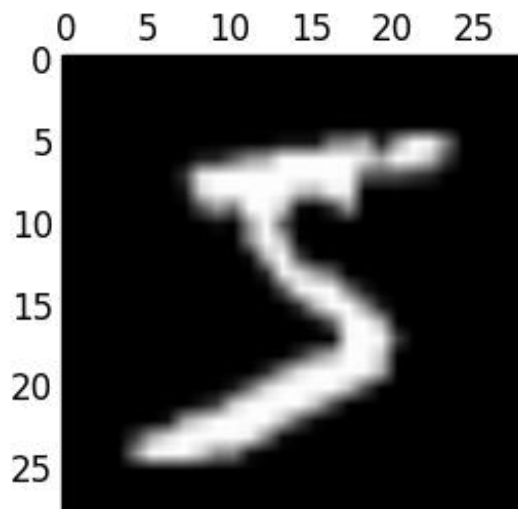


Modified National Institute of Standards and Technology (USA)

# MNIST

# (Lab) 20.py

- 60,000 training images + 10,000 testing images

- Input image : 28 * 28 pixels → 784 pixels

- 784 input dimension

- 10 classes (output: 0 ~ 9)

- Softmax

- 90.23% of recognition rate

https://github.com/yungbyun/neuralnetworks

Input Layer

Fully-connected

Output Layer

4 784

$\sum$ — $h_1$

$\sum$ — $h_2$

$\sum$ — $h_3$

$\sum$ — $h_4$

$\sum$ — $h_5$

Softmax

$\sum$ — $h_6$

$\sum$ — $h_7$

$\sum$ — $h_8$

$\sum$ — $h_9$

$\sum$ — $h_{10}$

10

Fully-connected

Softmax

10

784

# (Lab) 21.py

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

https://github.com/yungbyun/neuralnetworks

Fully-connected

784

256

256

Softmax

10

# (Lab) 22.py

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

- 97.23% of accuracy

https://github.com/yungbyun/neuralnetworks

# (Lab) 23.py



Fully-connected

784   512   512   512   512   10

- Applying initialization
- 6-layer deep neural networks
- 97.83% of accuracy

# Decision Boundry

deeper

wider

Hidden Layers
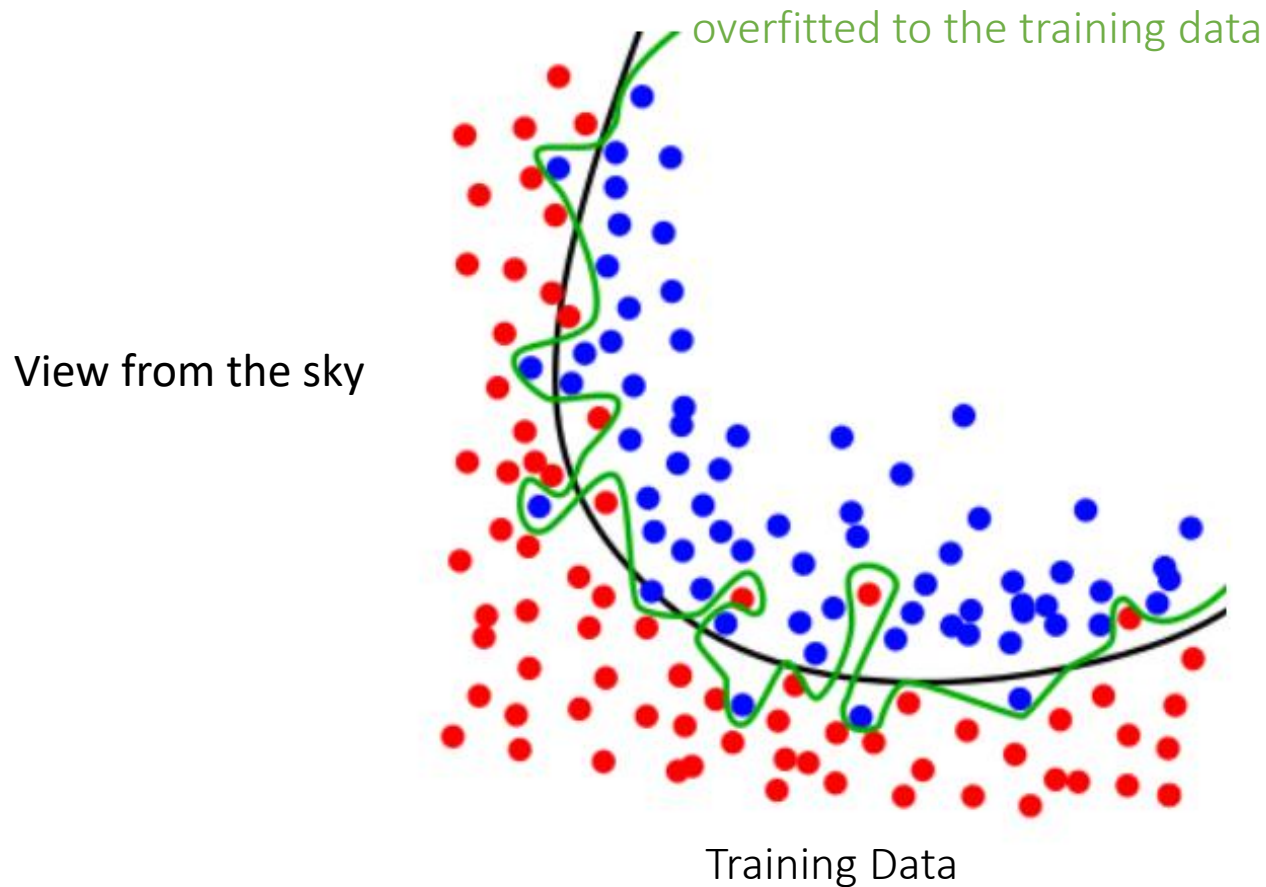
"

Lots of neurons (connections, synapses) give us
so complicated decision boundary.

# Which do you think is desirable decision boundary?

overfitted to the training data

View from the sky

Training Data

*While the black line fits the data well,*

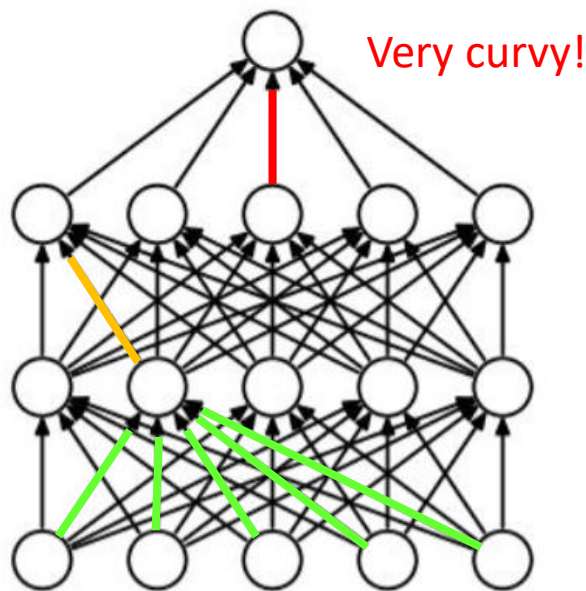*the green line is overfit.*

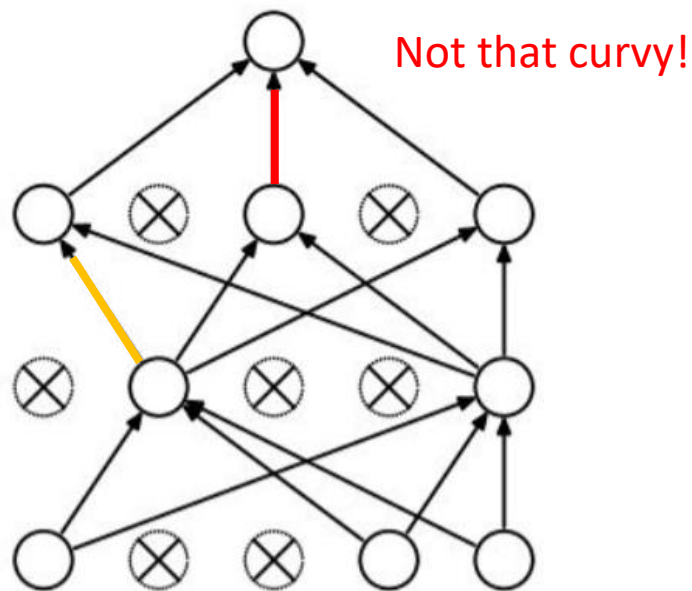https://elitedatascience.com

# Overfitting and drop-out

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

- Doing good at training data but errors for test data → overfitted to the training data

- Making it less complex by drop-out some neurons while learning.

# Regularization: **Dropout**

"randomly set some neurons to zero in the forward pass"



Very curvy!

Not that curvy!

(a) Standard Neural Net
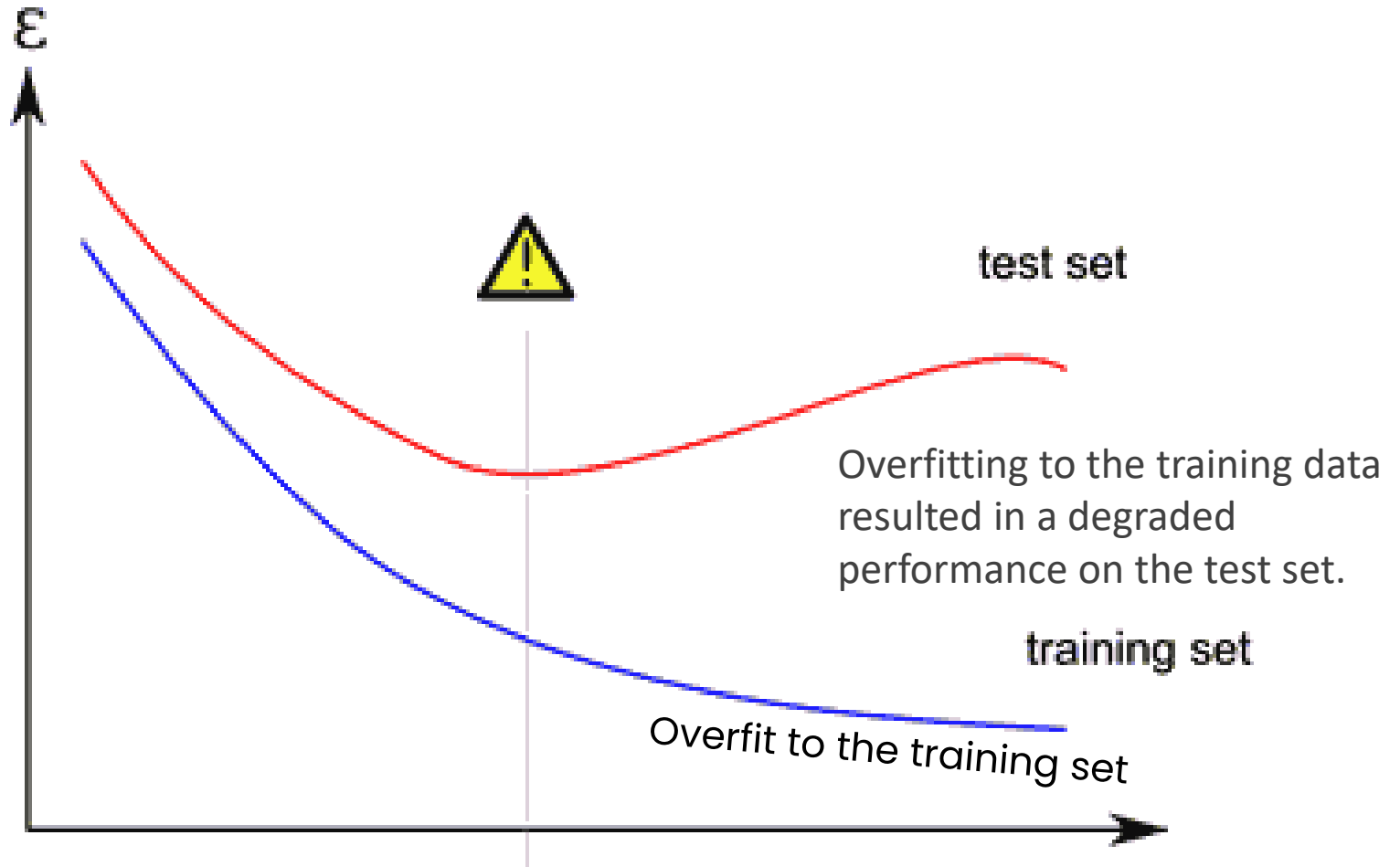
(b) After applying dropout.

[Srivastava et al., 2014]

# (Lab) 24.py

- Applying dropout
- 98.13% (←97.83%) of recognition accuracy

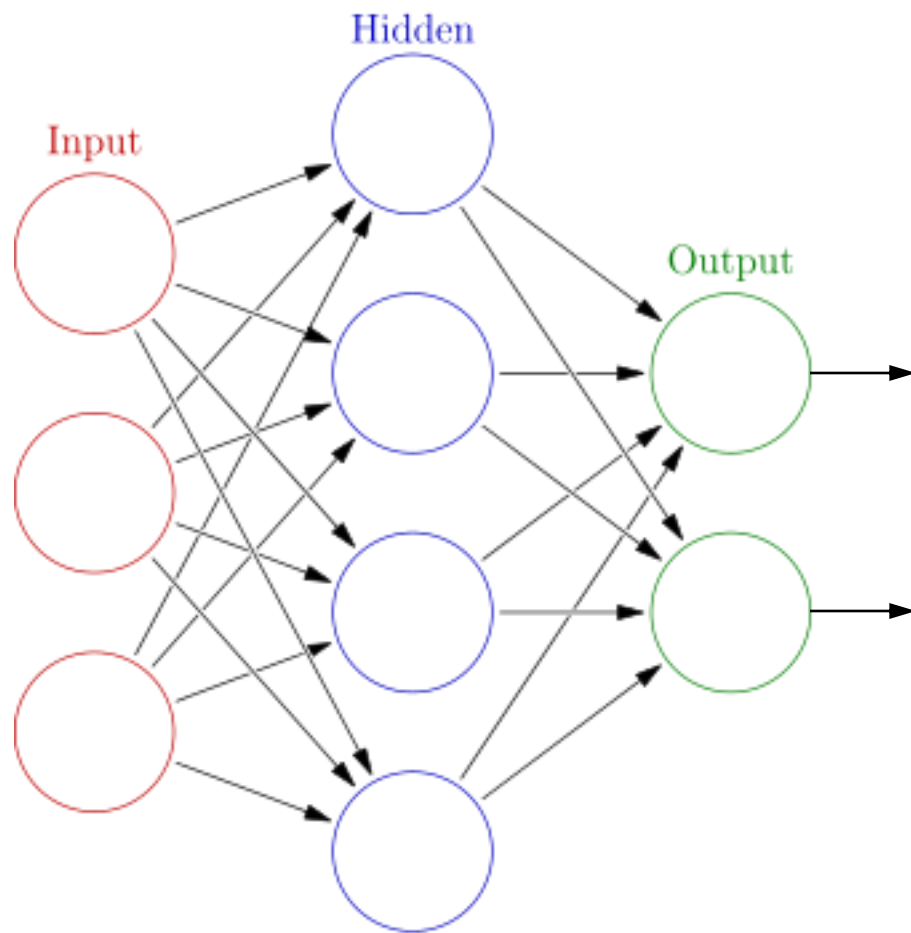# How to Prevent Overfitting

- Train with more data
- Reduce features
- Early stopping before overfitting
- Ensemble
- Regularization

# Early stopping



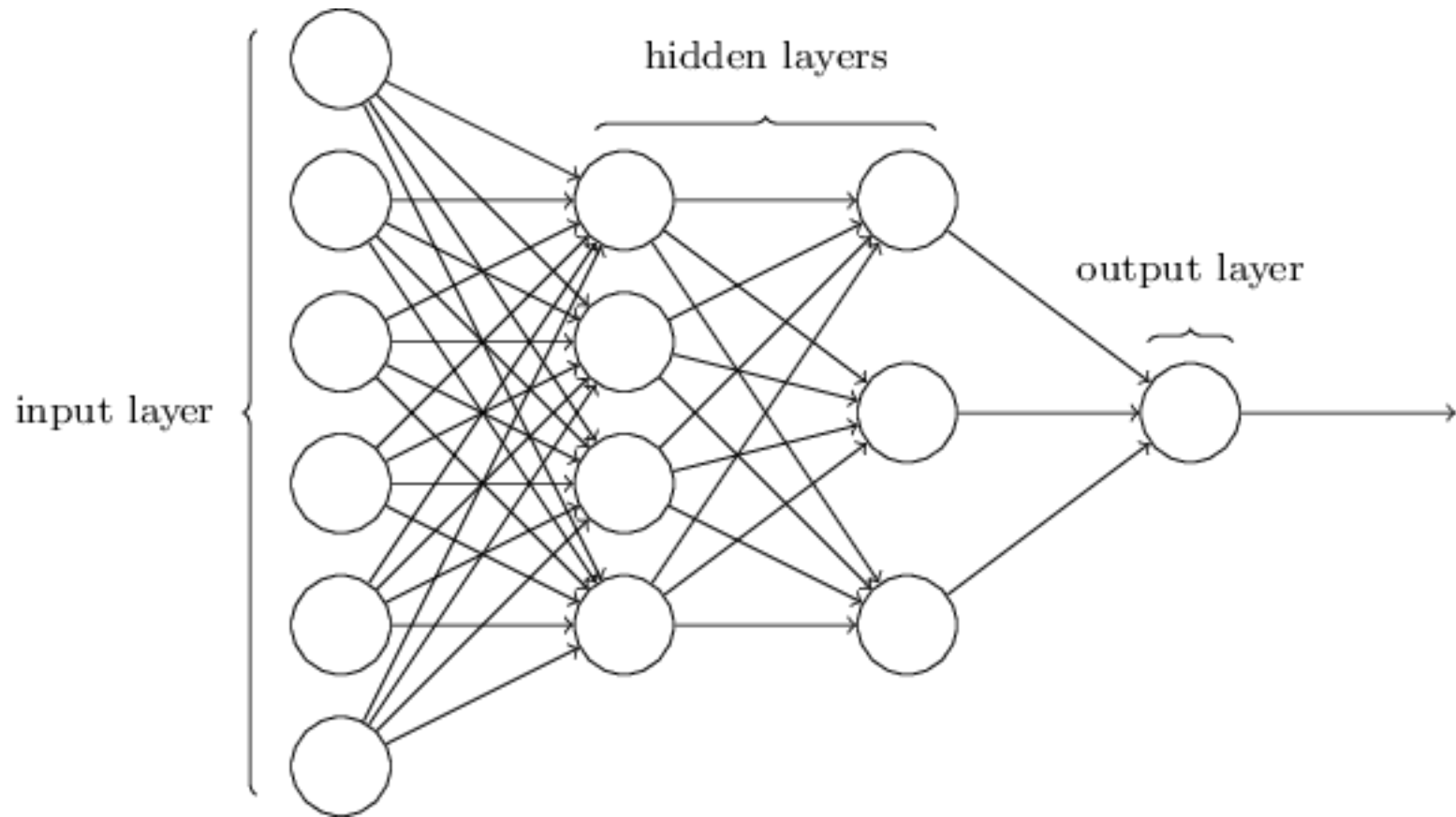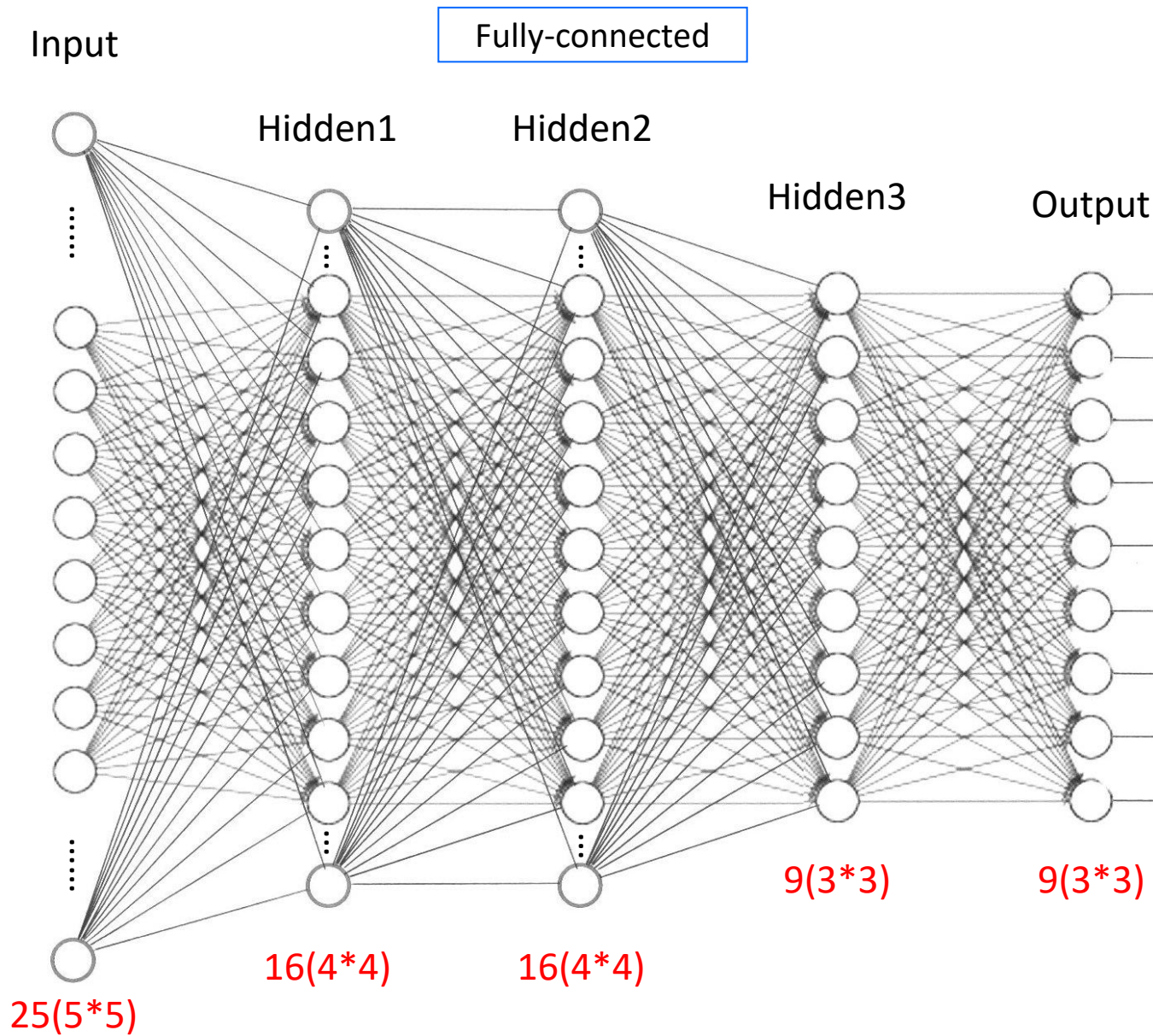$\varepsilon$

test set

Overfitting to the training data
resulted in a degraded
performance on the test set.

training set

Overfit to the training set

Fully-connected

Input

Hidden

Output

Fully-connected

hidden layers

output layer

input layer

Input

Hidden1　　Hidden2

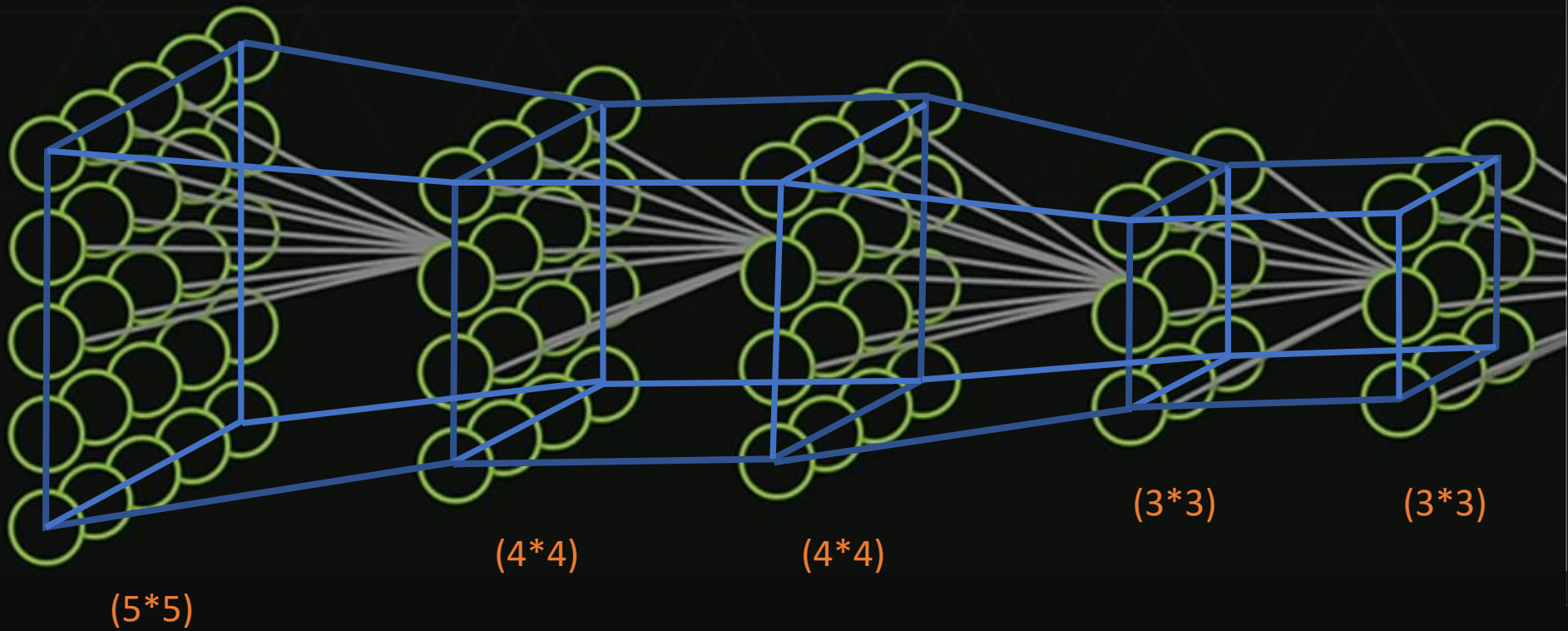Hidden3　　Output

Fully-connected

9(3*3)　　9(3*3)

16(4*4)　　16(4*4)

25(5*5)

Fully connected, then how many connections(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

(5*5)  (4*4)  (4*4)  (3*3)  (3*3)

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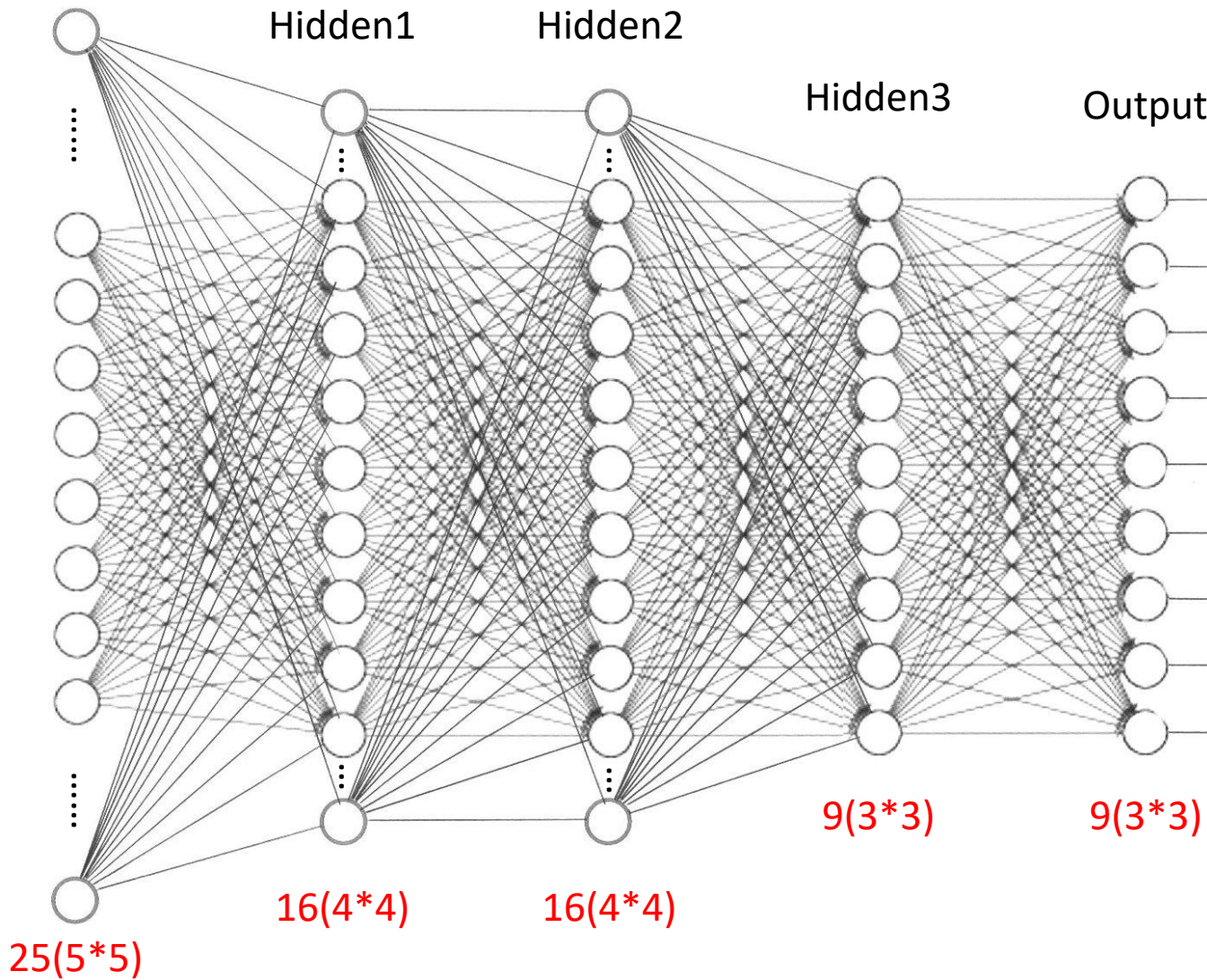25 * 16 + 16 * 16 + 16 * 9 + 9 * 9 = 881

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 connections(synapses, parameters) are there?
25 * 16 + 16 * 16 + 16 * 9 + 9 * 9 = 881

Geoffrey Hinton, Yann LeCun, Yoshua Bengio, Andrew Ng
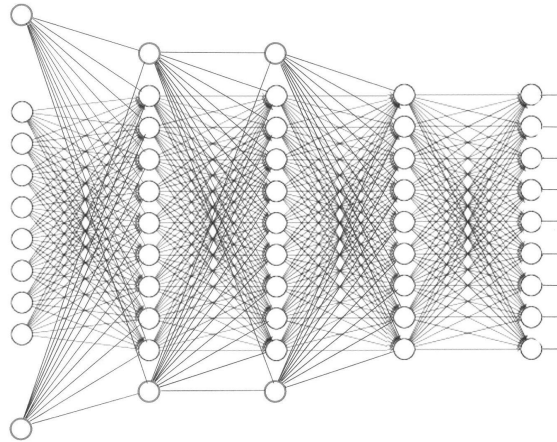
Google    Meta    Mila    amazon

# 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

Fully-connected Neural Networks

Any problem? Yes, so **CNN**.