

Perceptron & ANN

AI/DL areas and applications: past, current, upcoming



AI/DL Momentum

2016년 이세돌 vs. 알파고로 인해 대중들에게 AI 대한 관심도가 급상승 함.

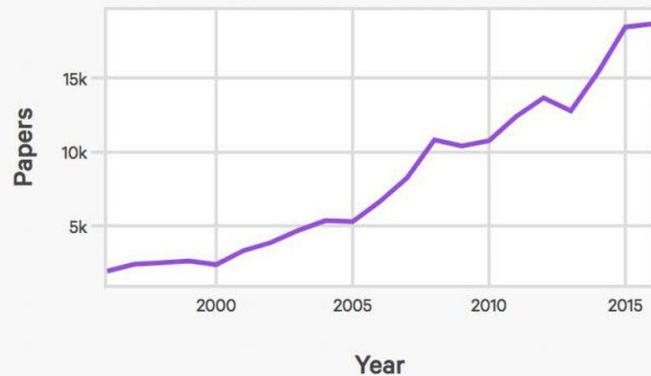


2022년 large language model (LLM) 기반 chaGPT가 등장함.



AI GROWTH

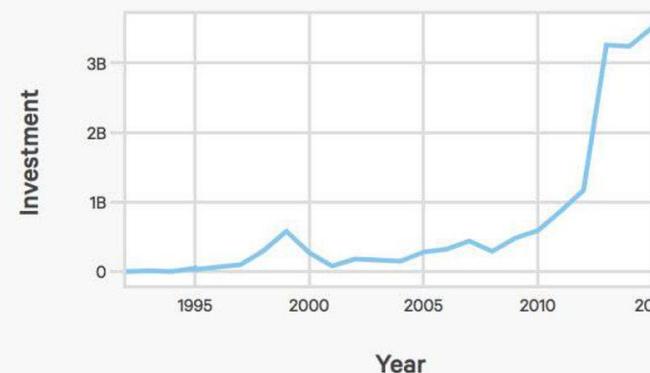
Annually Published AI Papers



Source: Scopus.com

AIINDEX.ORG

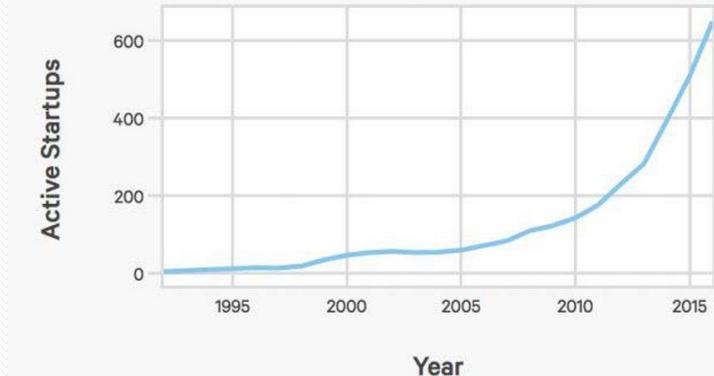
Annual VC Investment in AI Startups



Sources: Crunchbase, VentureSource, Sand Hill Econometrics

AIINDEX.ORG

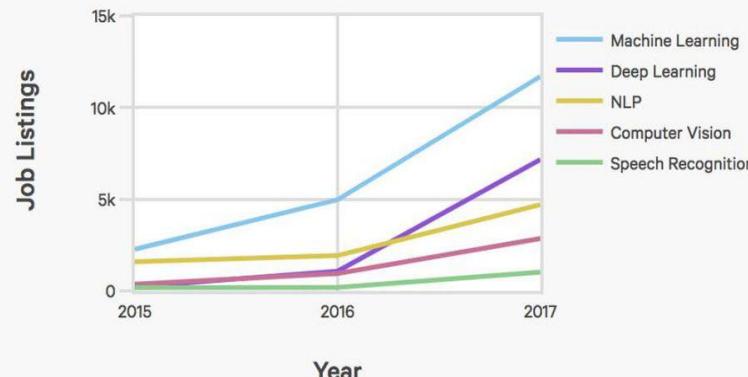
Startups Developing AI Systems



Sources: Crunchbase, VentureSource, Sand Hill Econometrics

AIINDEX.ORG

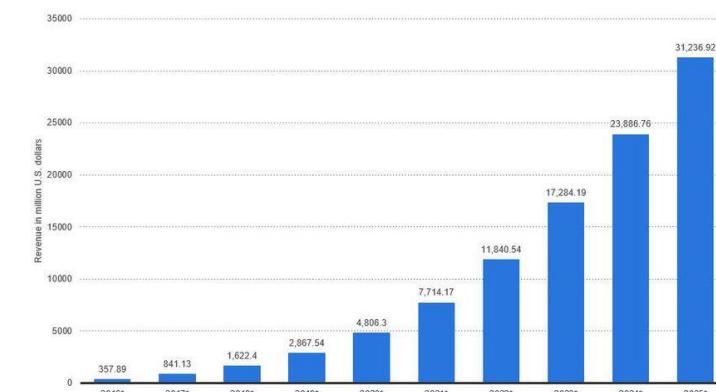
Job Openings, Skills Breakdown (Monster.com)



Source: Monster.com

AIINDEX.ORG

Enterprise artificial intelligence market revenue worldwide 2016-2025
Revenues from the artificial intelligence for enterprise applications market worldwide, from 2016 to 2025 (in million U.S. dollars)



statista

source

IMAGENET LARGE SCALE VISUAL RECOGNITION CHALLENGE ILSVRC

2010 - NEC-UIUC Lin et al.

2011 - XRCE Florent Perronnin, Jorge Sanchez

2012 - [AlexNet](#)

2013 - ZFNet

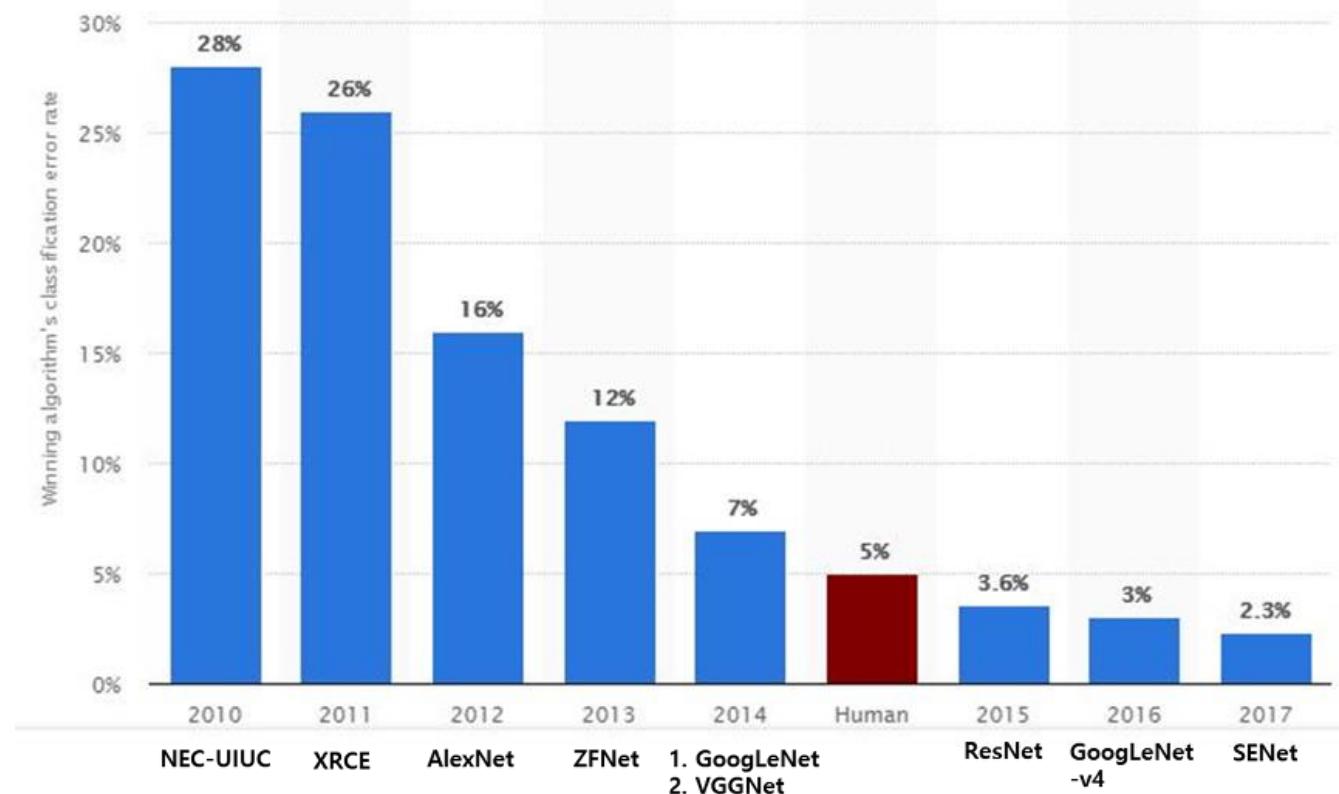
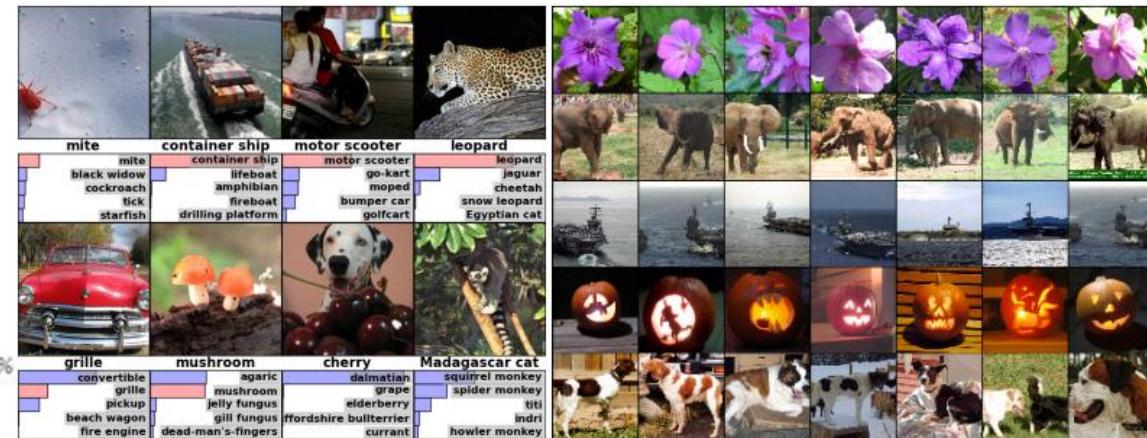
2014 – [GoogLeNet](#)

[VGGNet](#) Second Winner

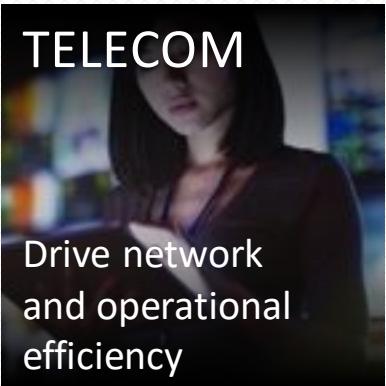
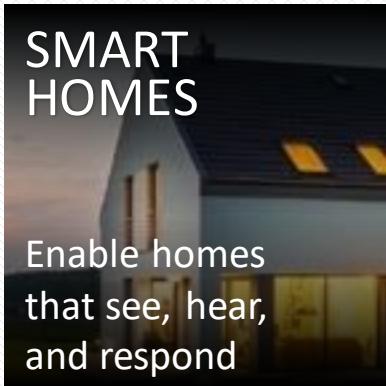
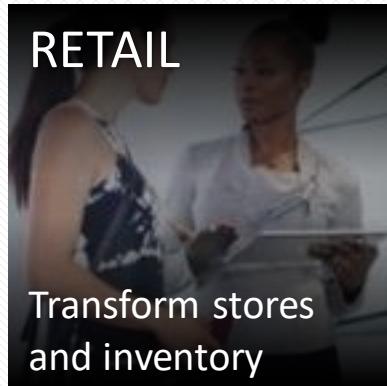
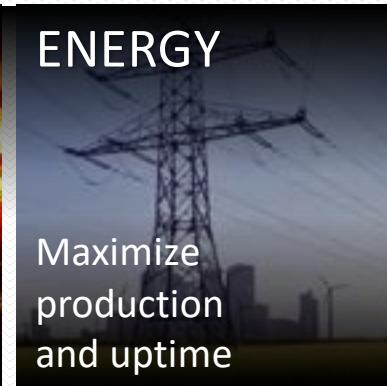
2015 - [ResNet](#)

2016 - GoogLeNet-v4

2017 - SENet



AI is changing every market

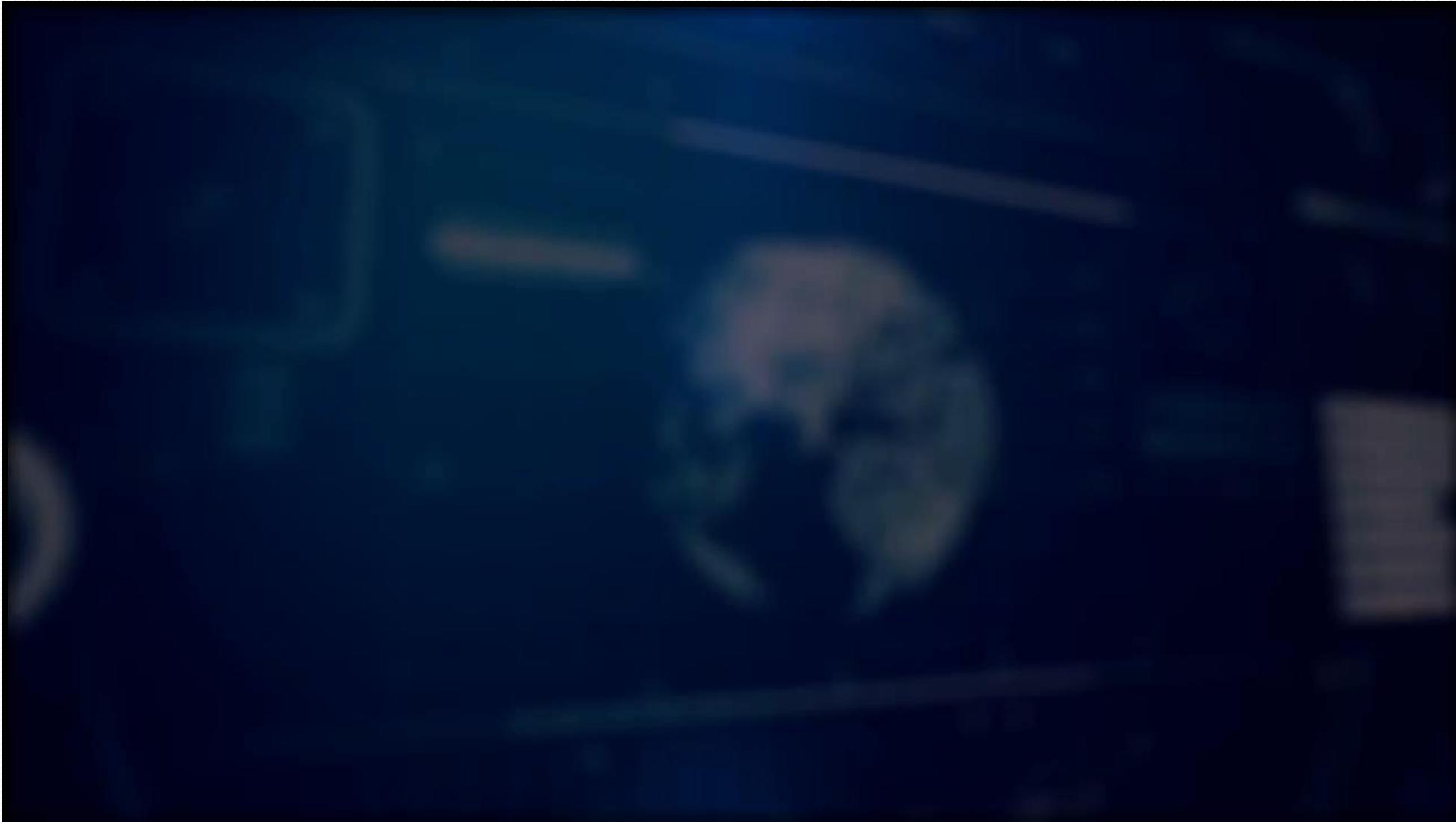
[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

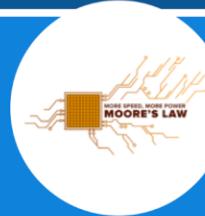
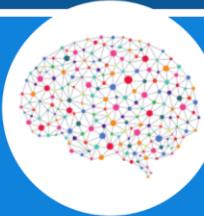
*Other names and brands may be claimed as the property of others.

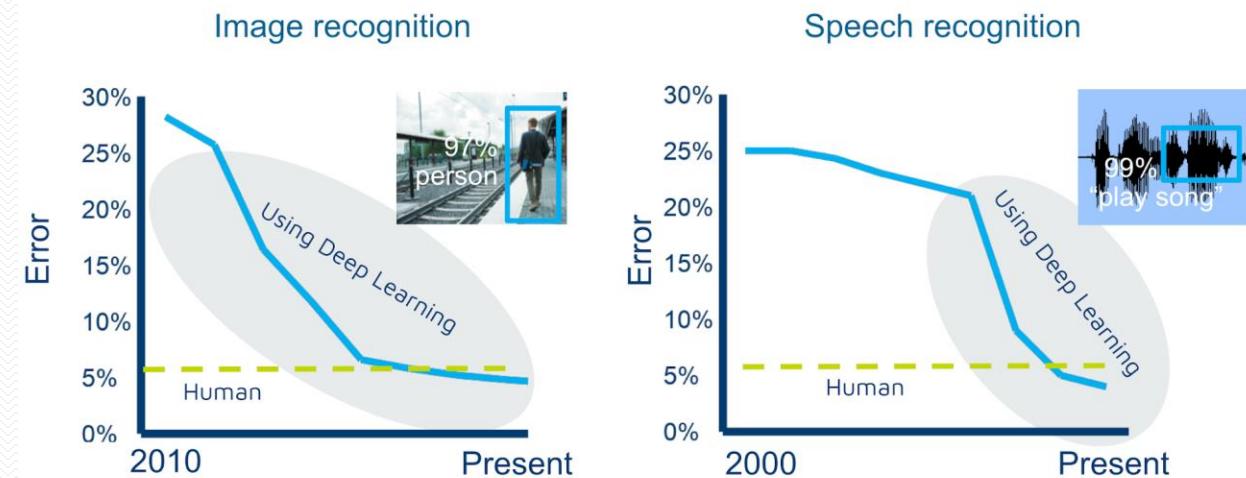
For public use – OK for non-NDA disclosure

WHY DO WE NEED AI?

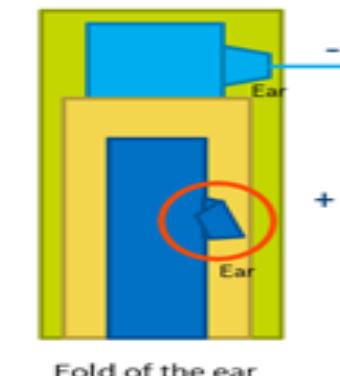
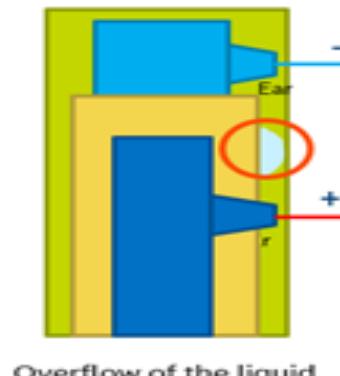
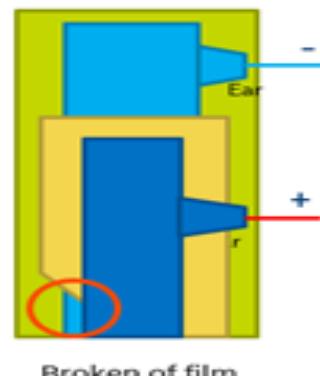
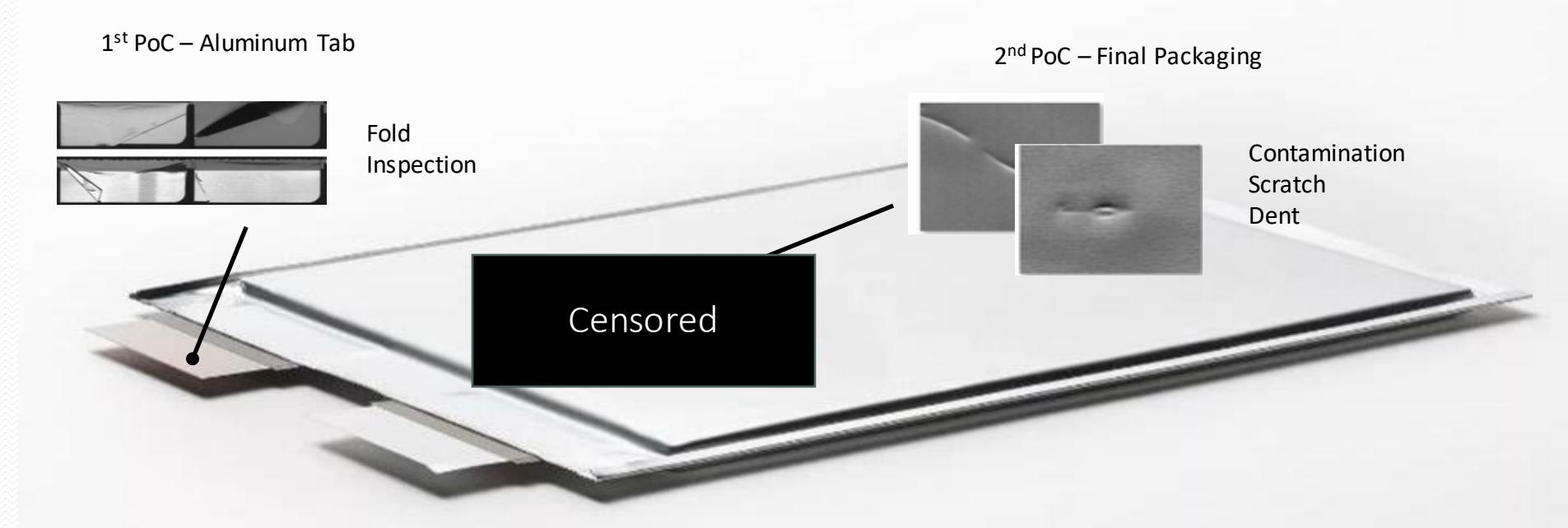


WHY NOW?

Bigger Datasets	Better Hardware	Smarter Algorithms
 IMAGENET 10M labeled images  Moore's Law Cost / GB in 1995: \$1000 Cost / GB in 2020: \$0.02	 Recurrent Neural Networks Convolutional Neural Nets LSTM	



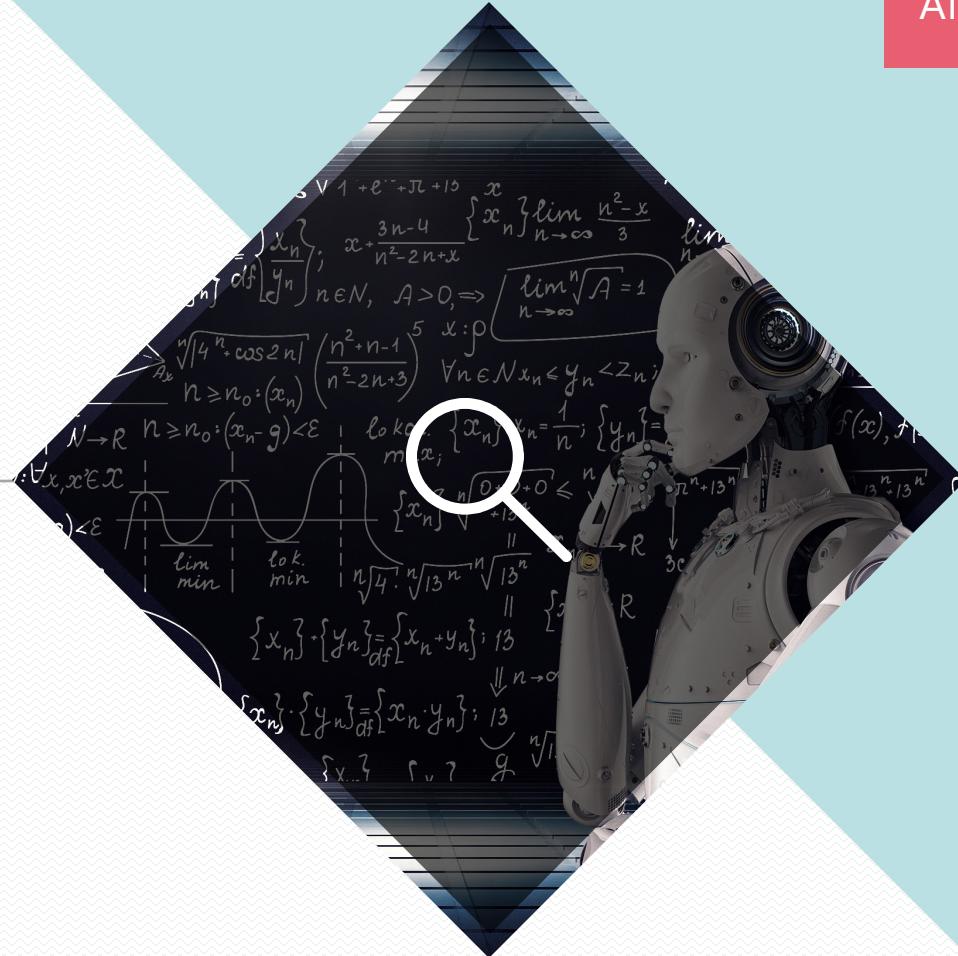
불량검출



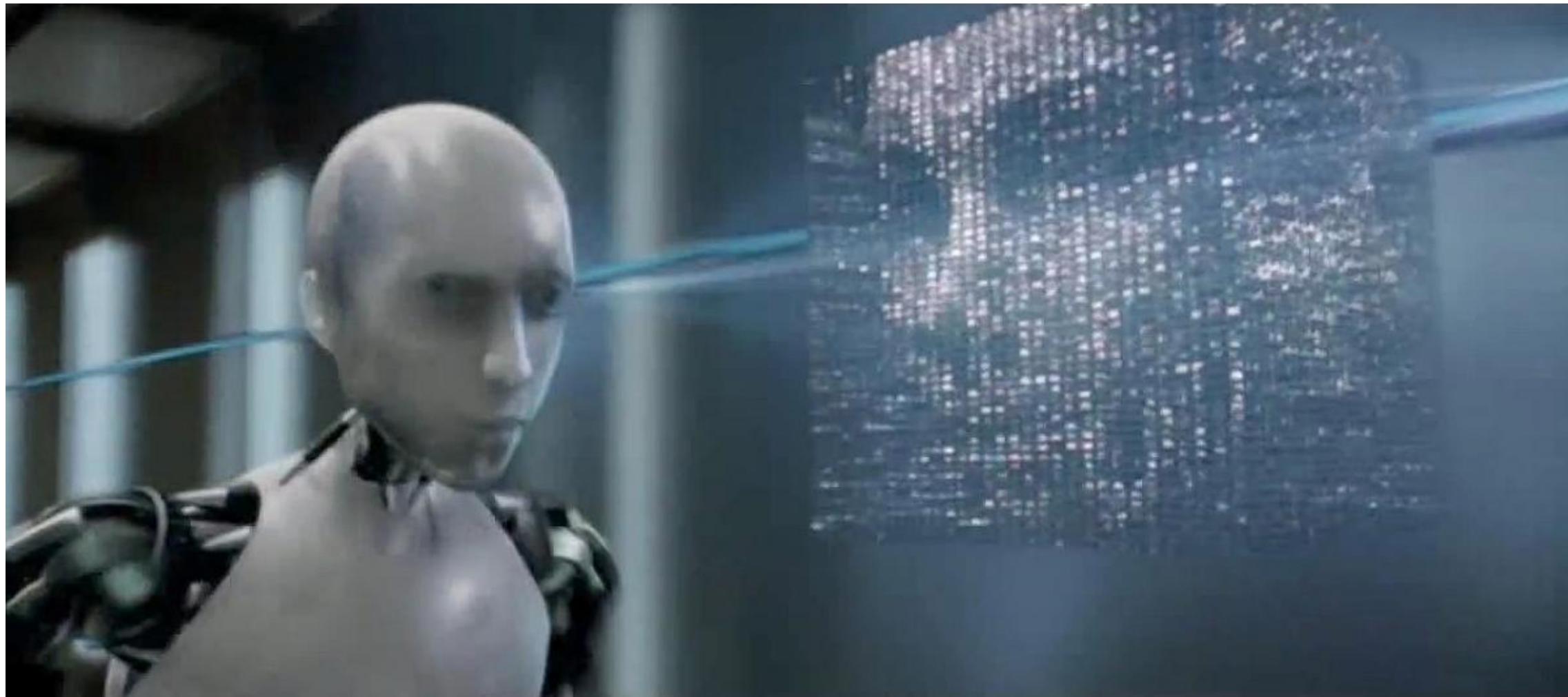
Session Break

AI

What is Artificial Intelligence? How Does AI Work?



AI 란 무엇일까?



자연현상에서 출발

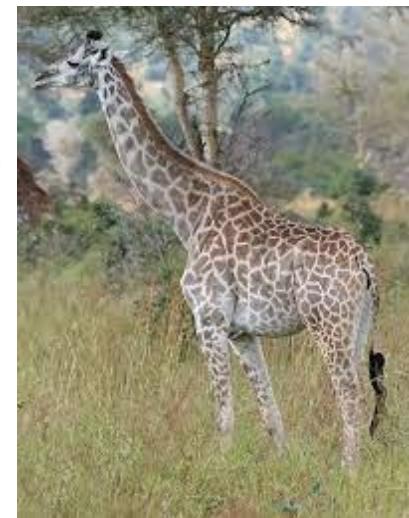
Sidewinder



Ultrasound



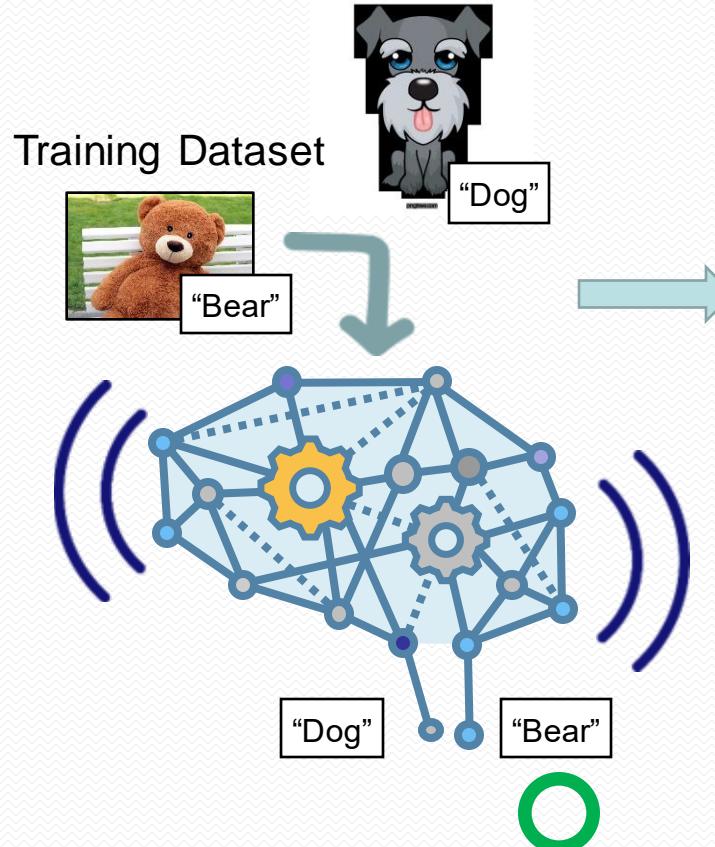
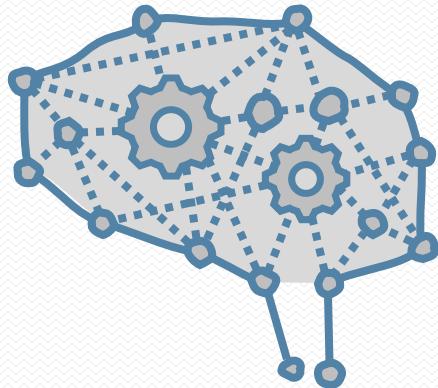
우리는 어떻게 생각했나



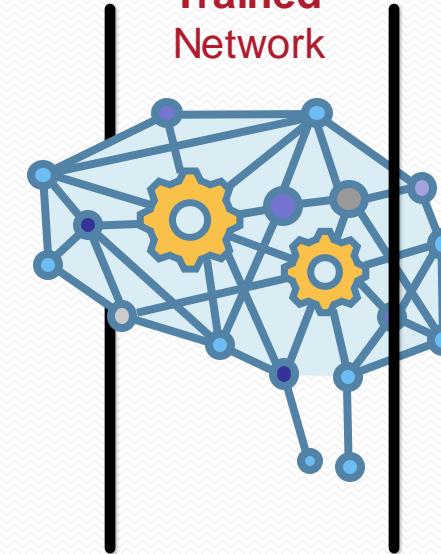
Deep Learning

TRAINING

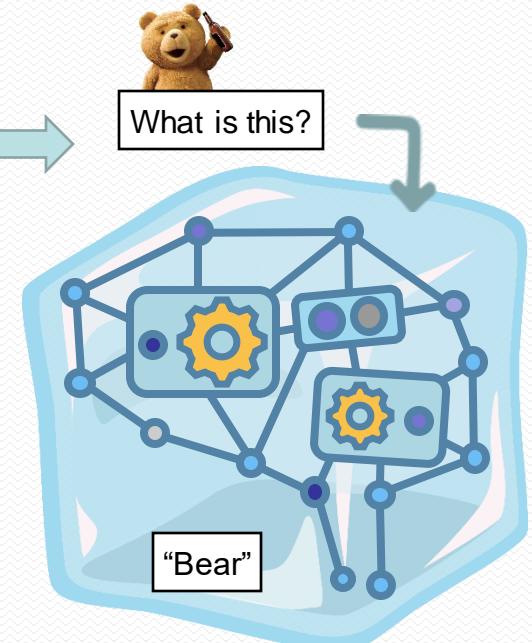
Untrained Network



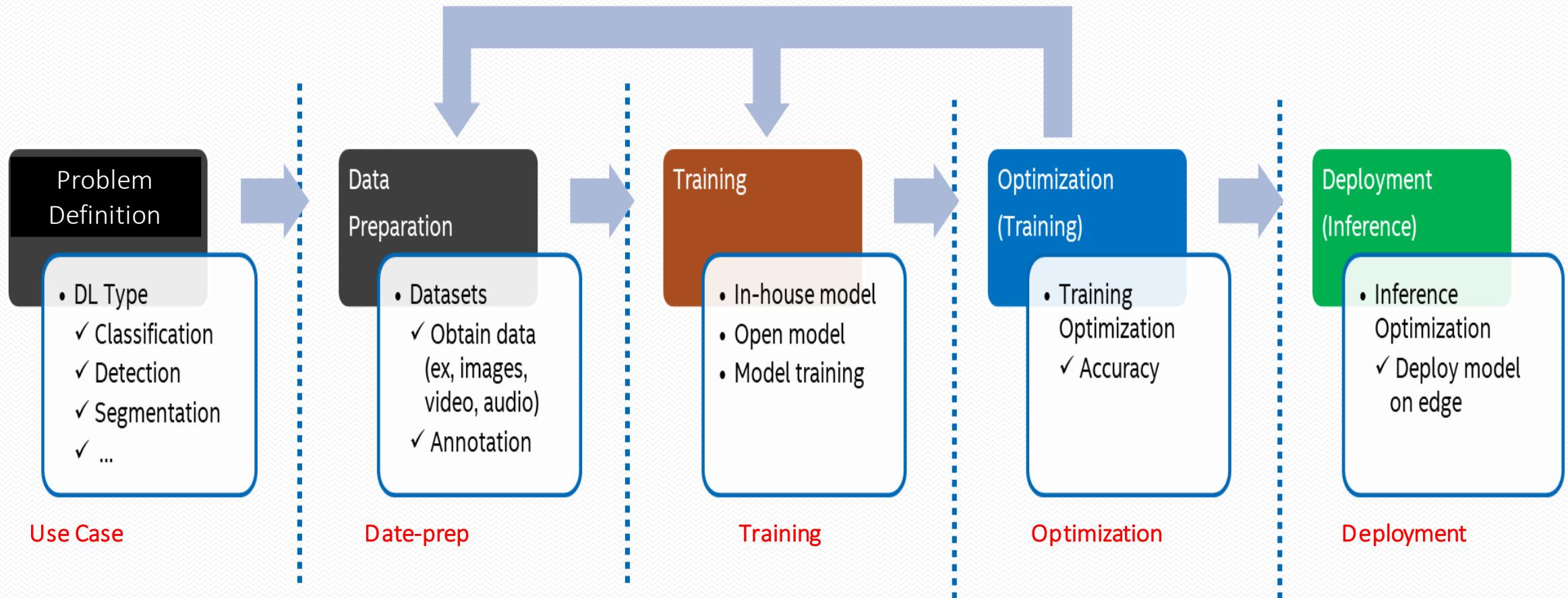
Trained Network



INFERENCE



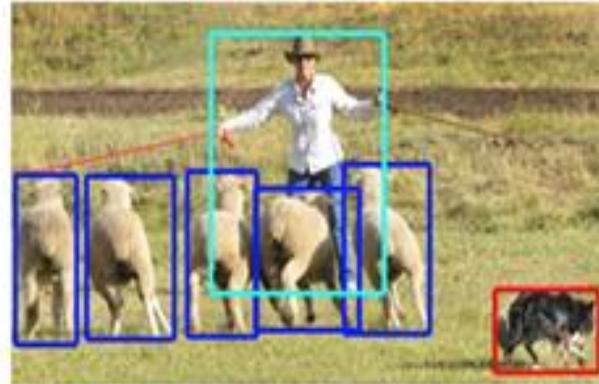
Deep Learning Workflow (생각하게 만들기)



1) Problem Definition



CLASSIFICATION



DETECTION



SEGMENTATION



NEURAL STYLE TRANSFER



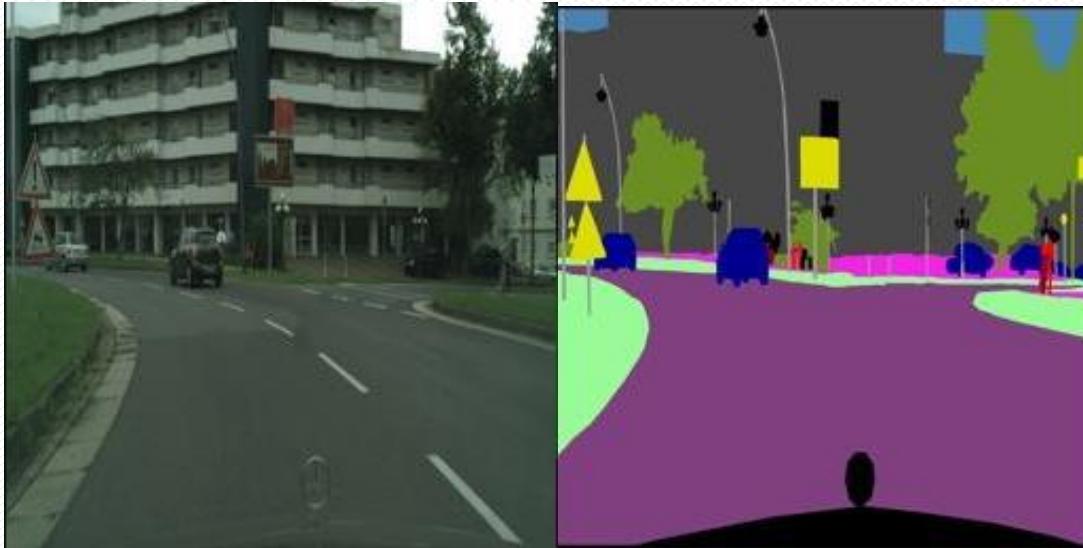
IMAGE CAPTION



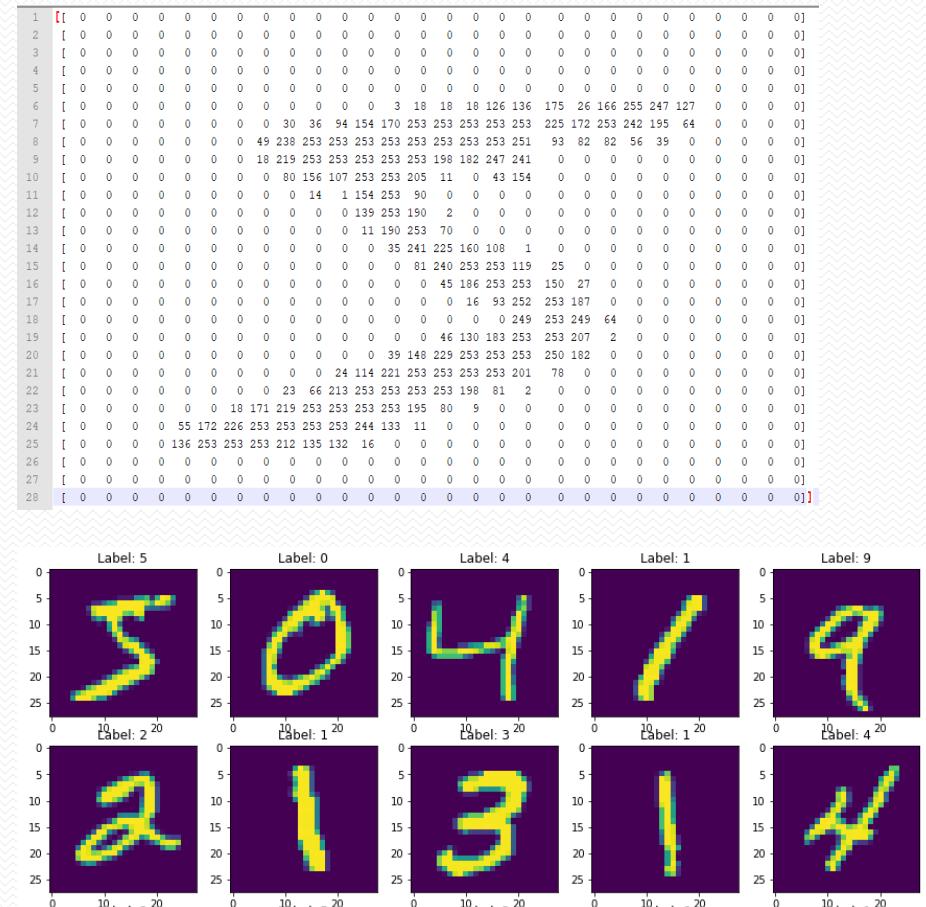
COLORIZATION

2) Dataset preparation/annotation

- Cityscape/Camvid
Semantic Segmentation



- MNIST
Hand-written digit



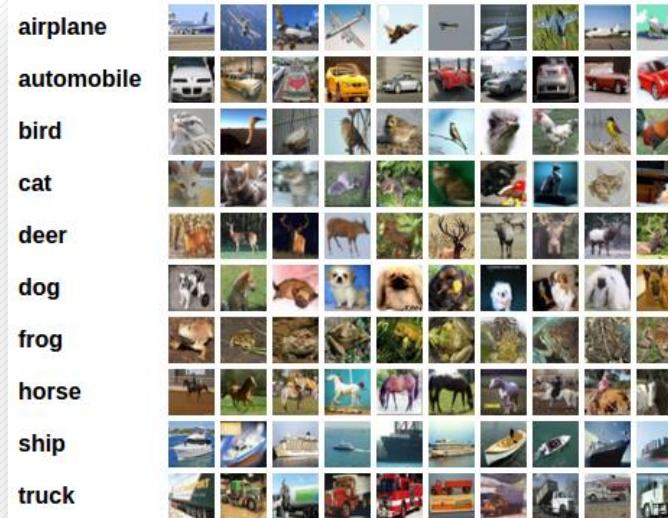
2) Dataset preparation/annotation – cont'd

- MNIST
 - Label : digit 0-9
 - Consists of 60,000 training image & 10,000 test image 28x28 gray scale image

0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

Data & Labels

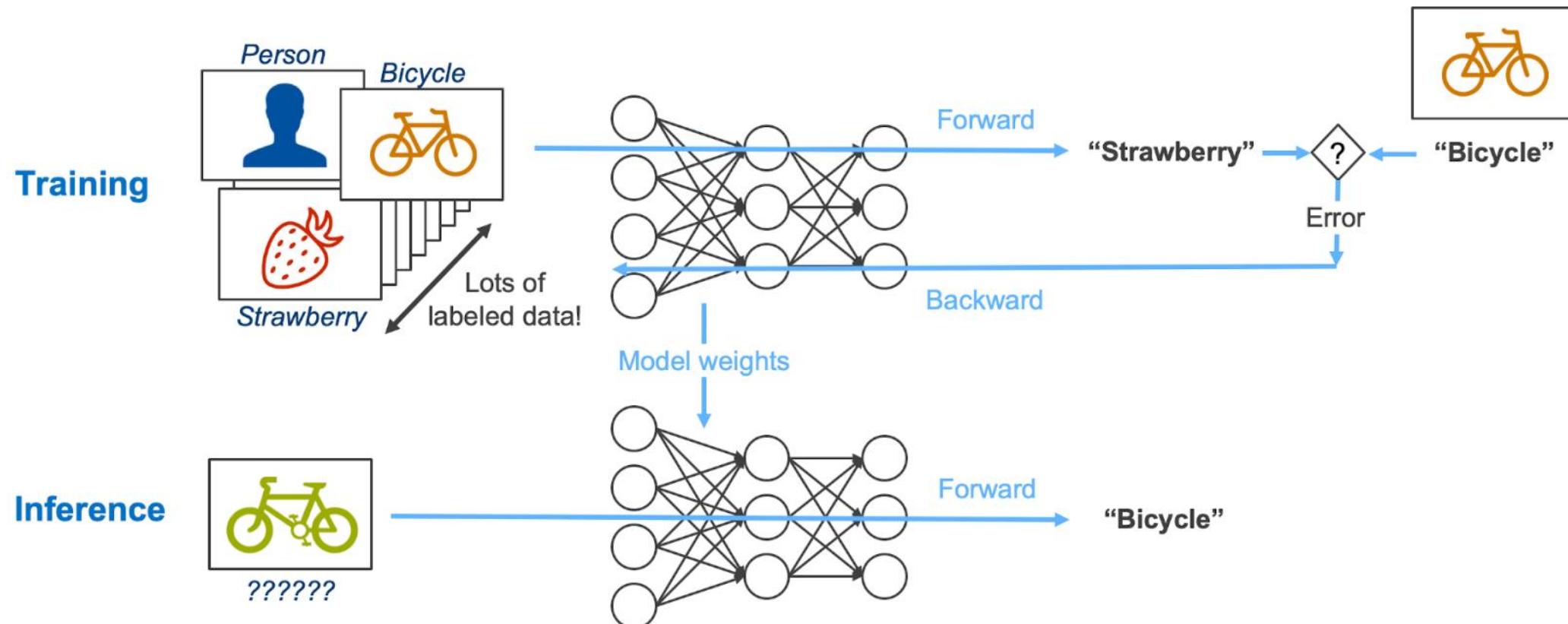
- CIFAR-10 dataset
 - Label : 10 classes
 - Consists of 60,000 (32x32) color images, with 6000 images per class.



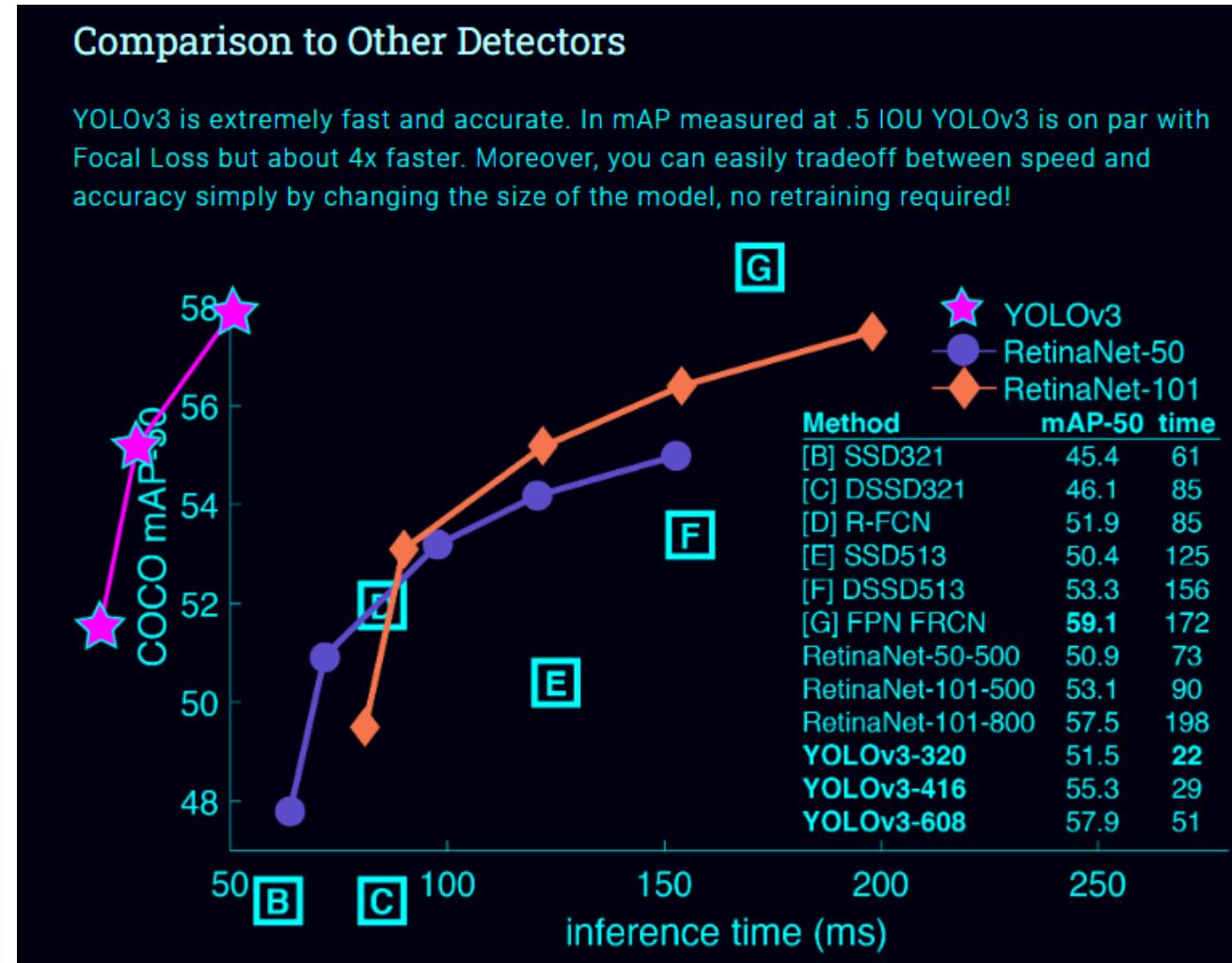
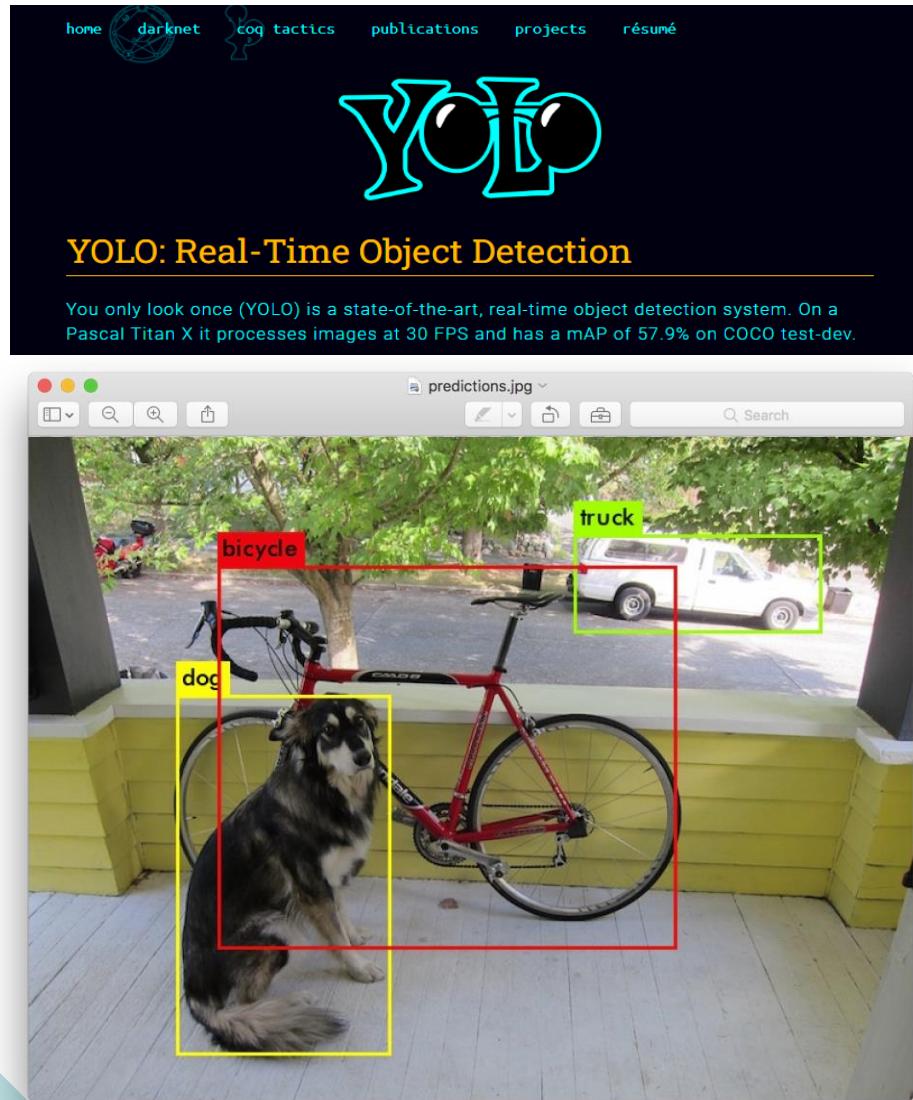
- IMAGENET
 - Label : 1000 classes
 - Consists of 800,000 (256x256 or 224x224) color Images, with 800 images per class.



3) Model selection & Training



Deep Learning Model (YOLO)



4) Optimization

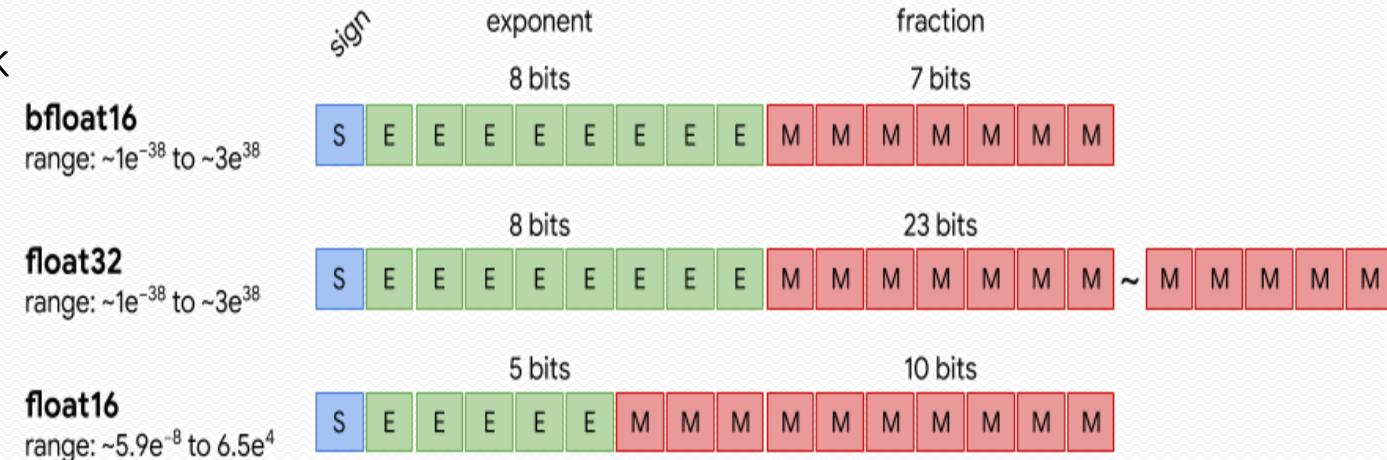
FP16 VS. FP32

FP16 improves speed (TFLOPS) and performance

FP16 reduces memory usage of a neural network

FP16 data transfers are faster than FP32

But converting from/to 32-bit floating-point can reduce accuracy



INT8 QUANTIZATION

Accelerate the performance of certain models

However, this comes at the cost of a small reduction in accuracy

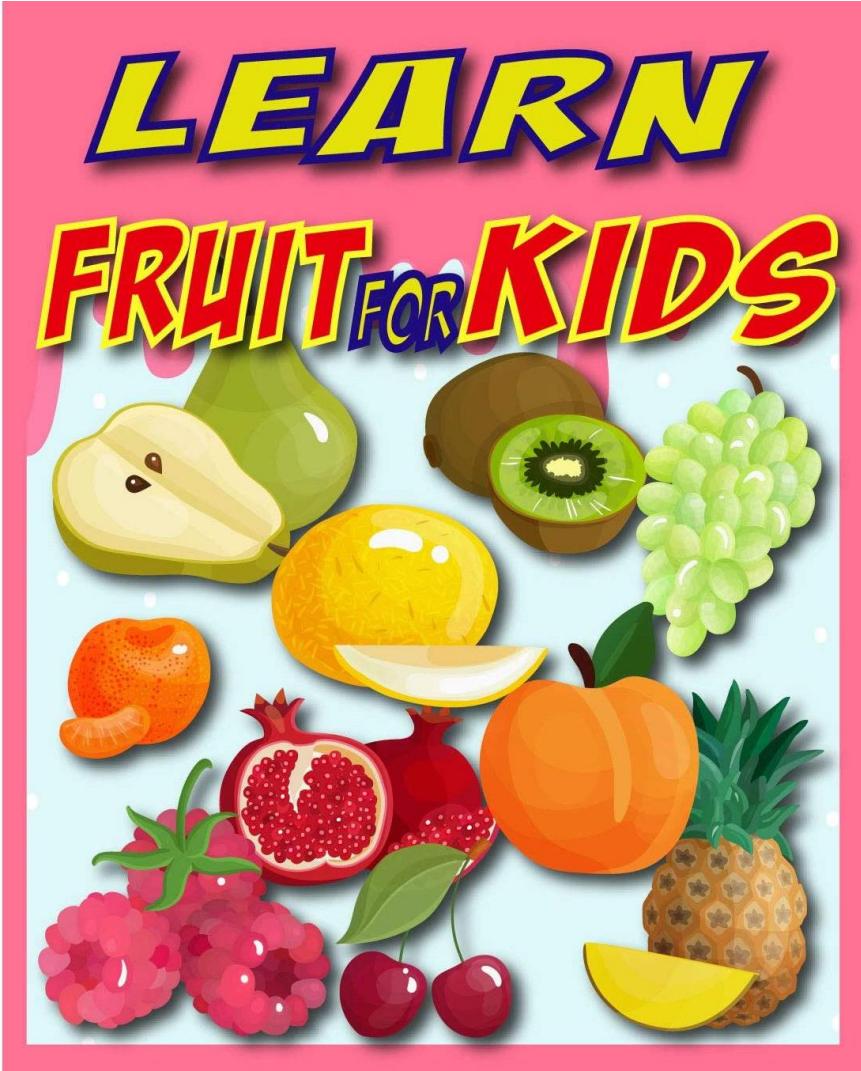
4) Optimization



사과!

??

4) Optimization

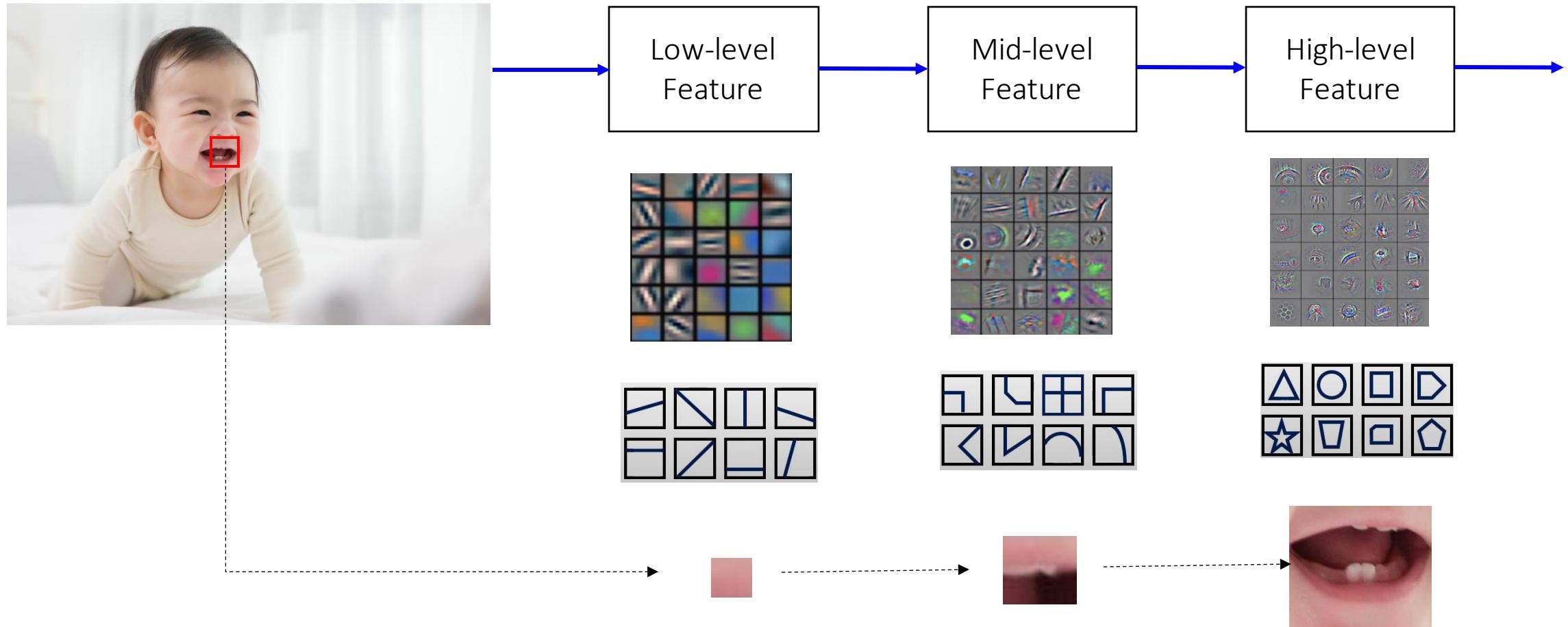


Apple!



Pear!

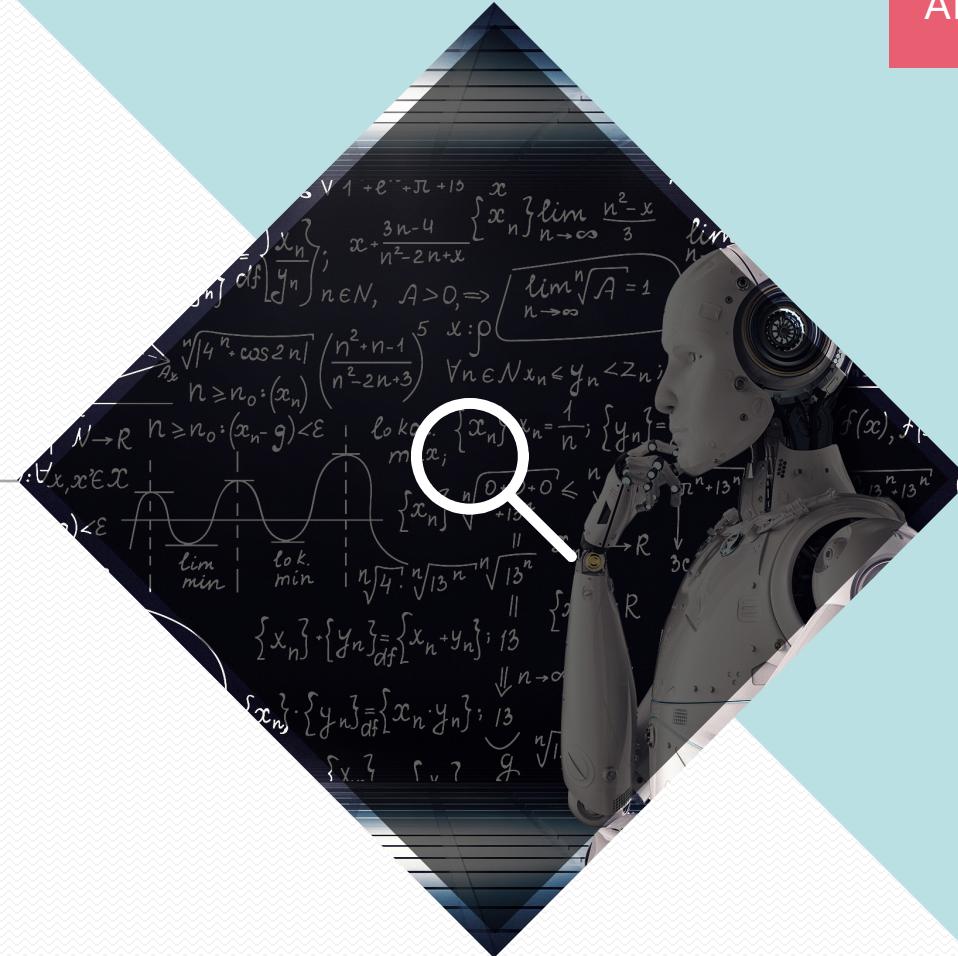
4) Optimization : transfer learning



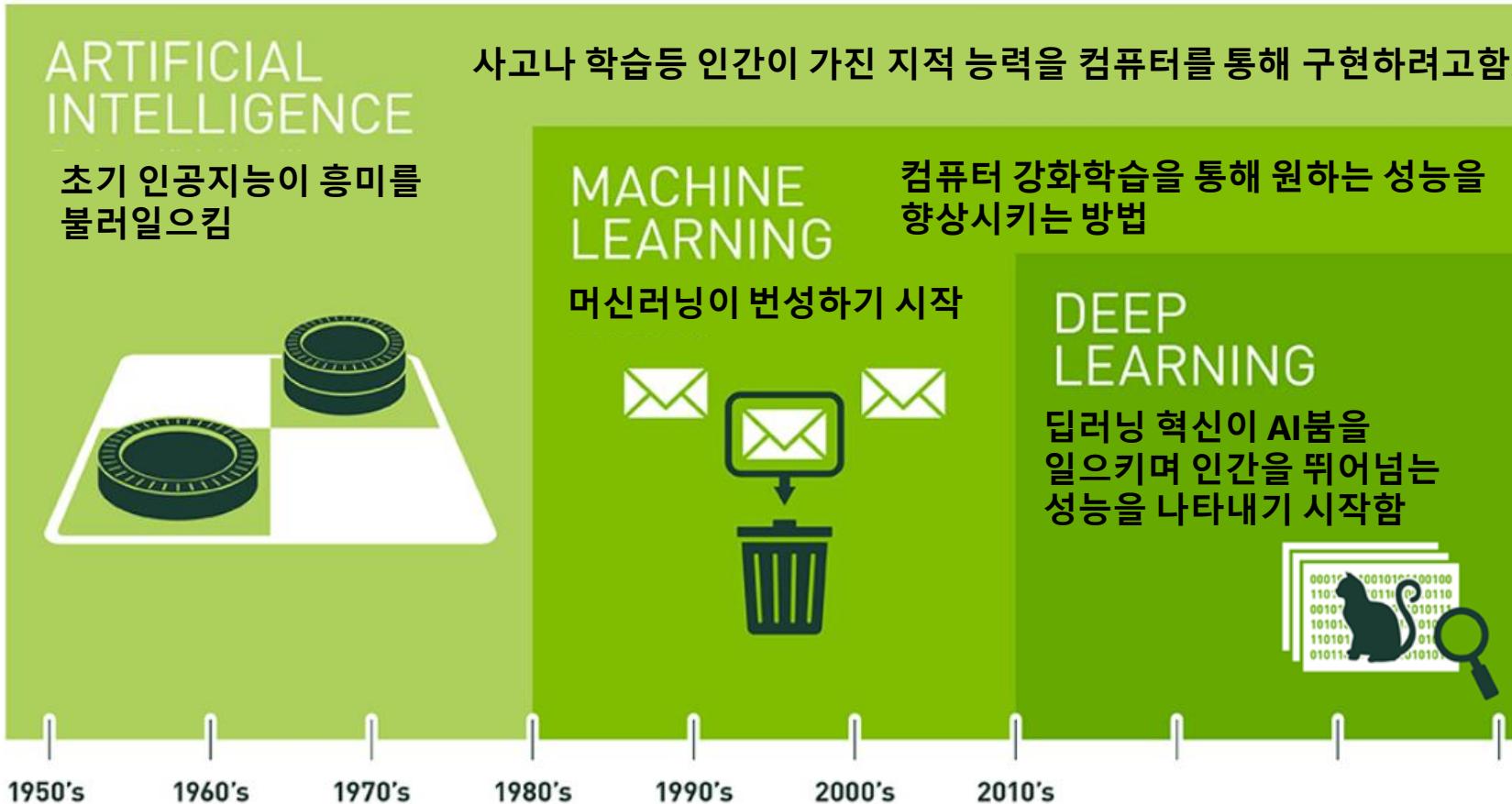
Session Break

History

How DL and AI have evolved?



인공지능



a system that is
“intelligent” through rules

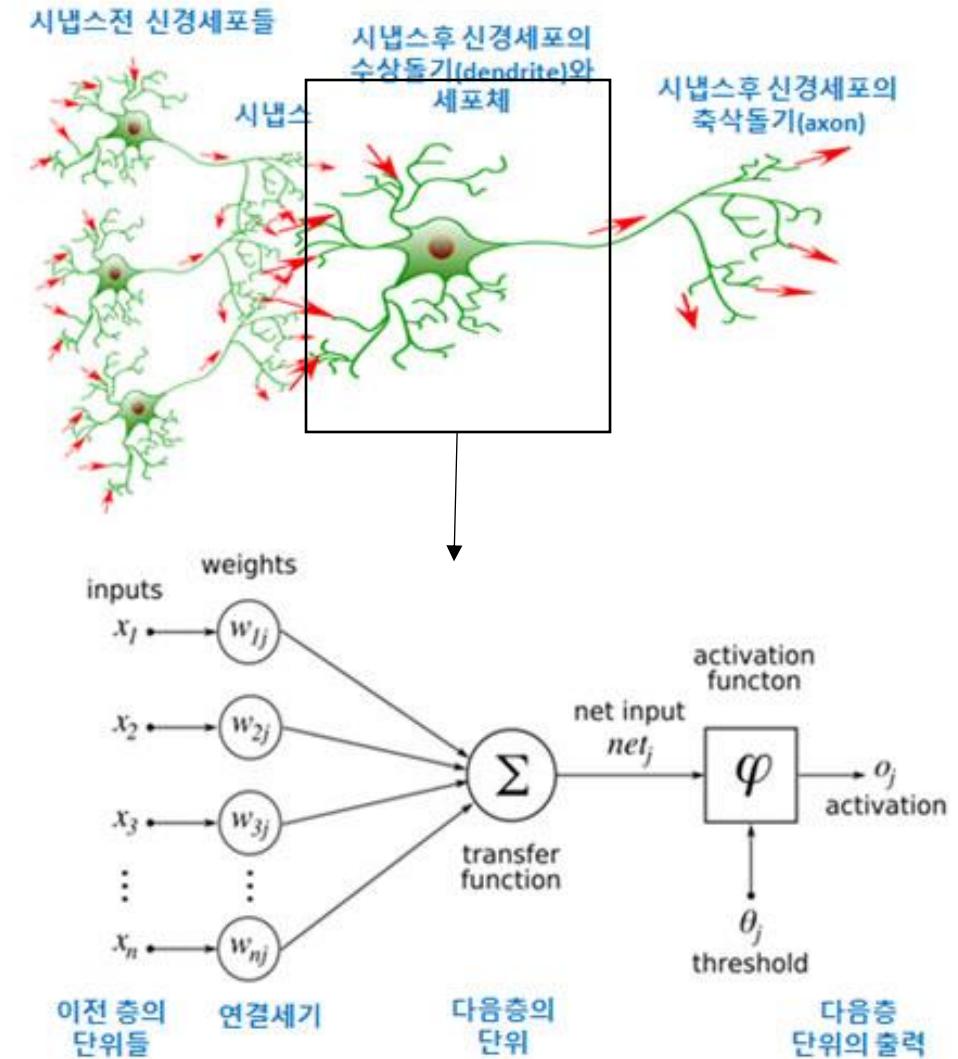
self-learning algorithms that
learn models from “data”

인공 신경망



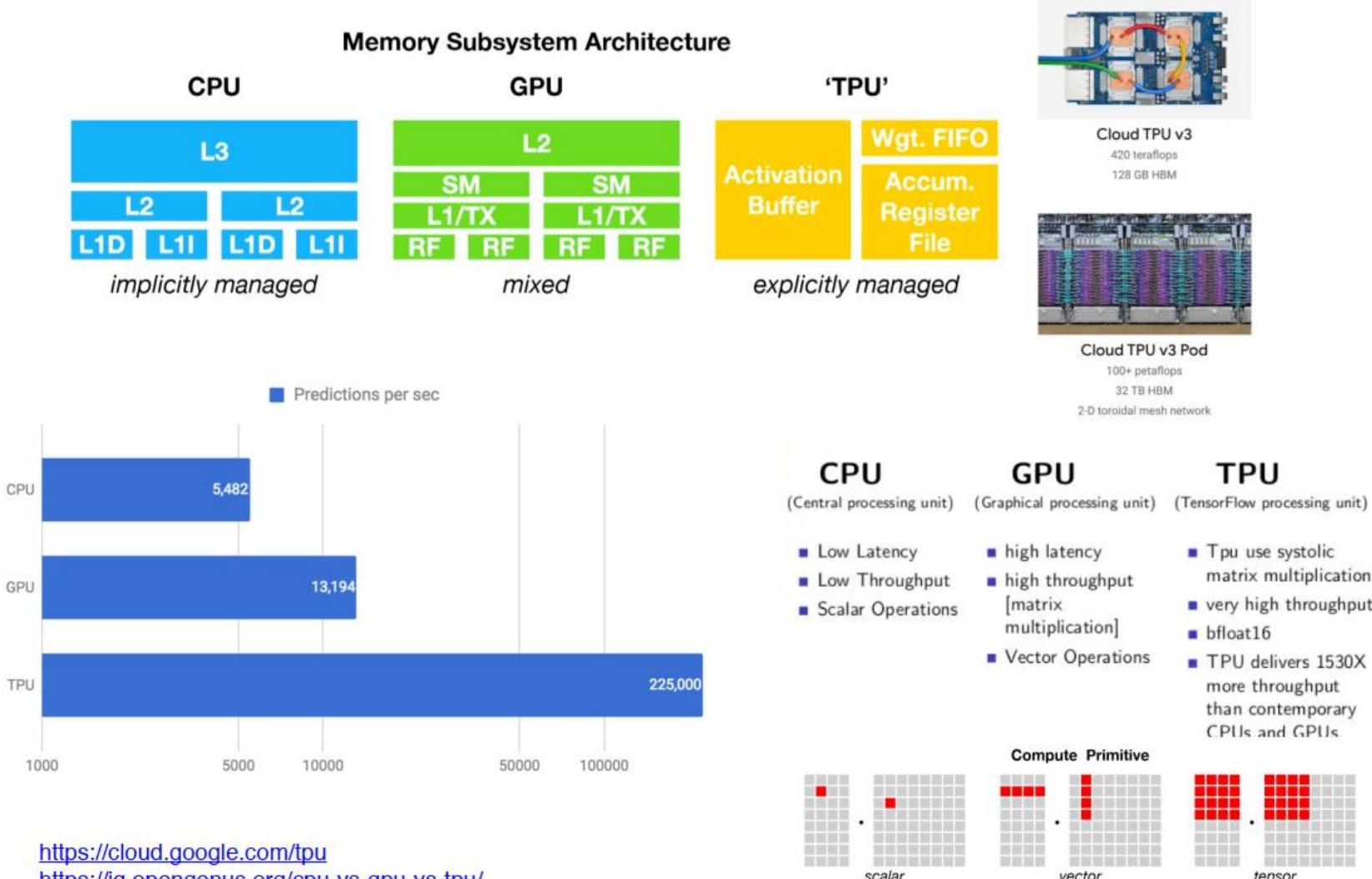
인공 신경망

뇌 속 신경망



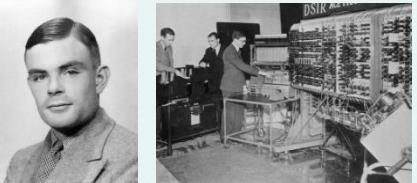
AI computing

CPU vs GPU vs TPU (Google Tensor Processing Unit)



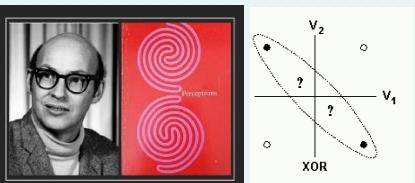
HISTORY OF A.I.

Alan Turing
Turing Test



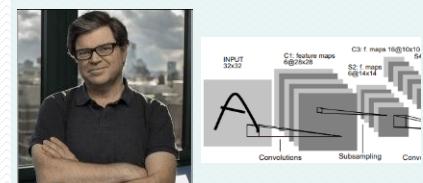
1950

Marvin Minsky
A book of "Perceptrons"



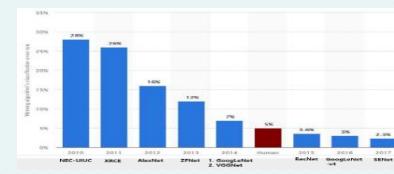
1969

Yann Lecun
LeNet-5



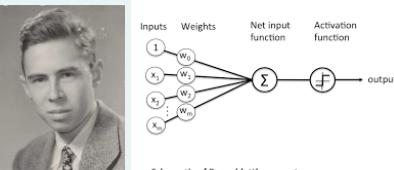
1998

ResNet
Accuracy is higher than human



2015

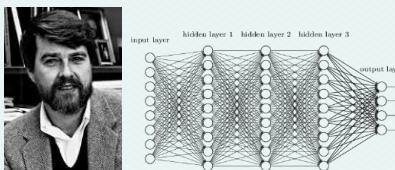
Frank Rosenblatt
Perceptron



1958

AI Winter

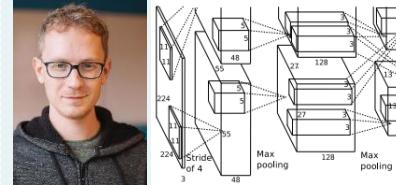
David Rumelhart
Multi Layer Perceptron (MLP)



1986

AI Winter

Alex Krizhevsky
AlexNet



2012

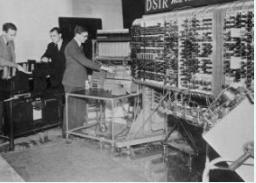
Google Brain
AlphaGo



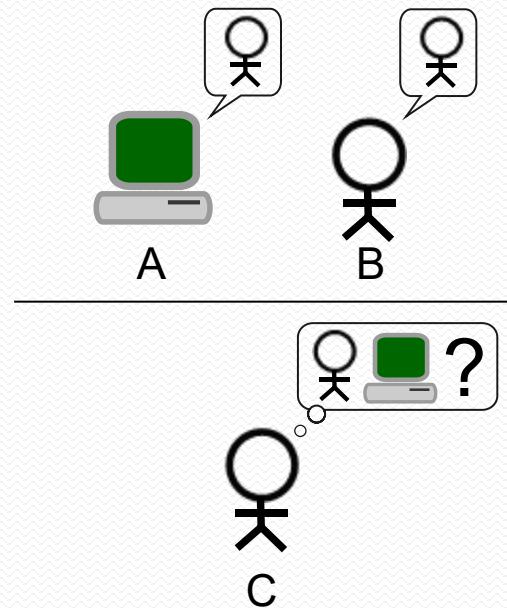
2016

HISTORY OF A.I.

Alan Turing
Turing Test



1950

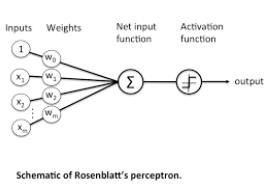


- 튜링 테스트(Turing test)의 의미
 - 기계가 생각하는 것이 가능한가?
 - Computing Machinery and Intelligence 논문
- 각각 분리된 채, 하나는 인간, 다른 하나는 기계와 대화를 하는 Test
- 5분동안 기계와 대화한다는 것을 인지하지 못하면 성공
- 대화에서 질문에 올바른 대답을 하는것 보단 얼마나 인간같이 대화할 수 있느냐로 판단

HISTORY OF A.I.

1958

Frank Rosenblatt
Perceptron



Schematic of Rosenblatt's perceptron.

- 프랭크 로젠블렛(신경 생물 과학자, 심리학전공)
- 이미지 인식 성공 !! – 많은 투자와 관심
- Feed Forward Network
- 발표 논문
 - *THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN*

If we are eventually to understand the capability of higher organisms for perceptual recognition, generalization, recall, and thinking, we must first have answers to three fundamental question.

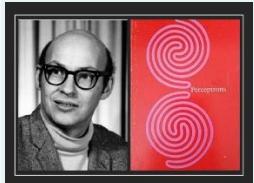
1. *How is information about the physical world sensed, or detected, by the biological system?*
2. *In what form is information stored, or remembered?*
3. *How does information contained in storage, or in memory, influence recognition and behavior?*



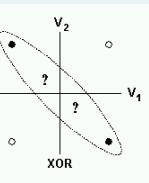
HISTORY OF A.I.

Marvin Minsky

A book of "Perceptrons"

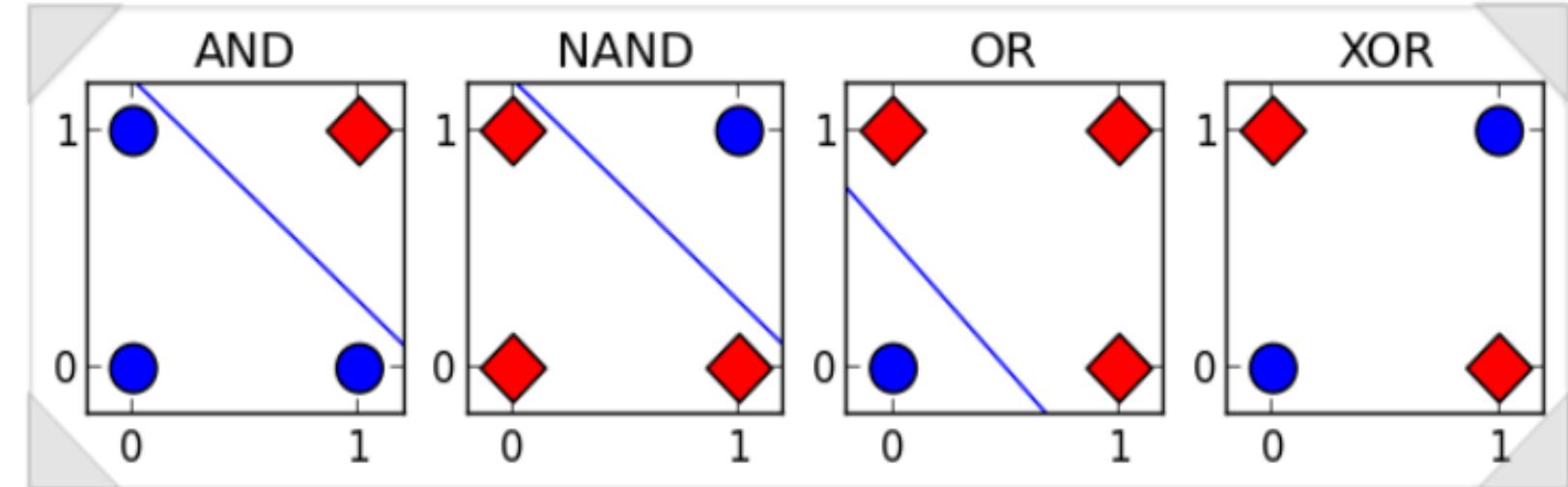


1969



마빈 민스키

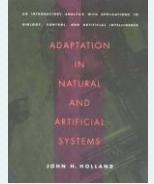
- *Perceptrons* – Minsky & Papert 출간, 퍼셉트론 모델에 대해 철저한 분석 및 그 한계에 대해서도 논리정연하게 분석



HISTORY OF A.I.

John Holland

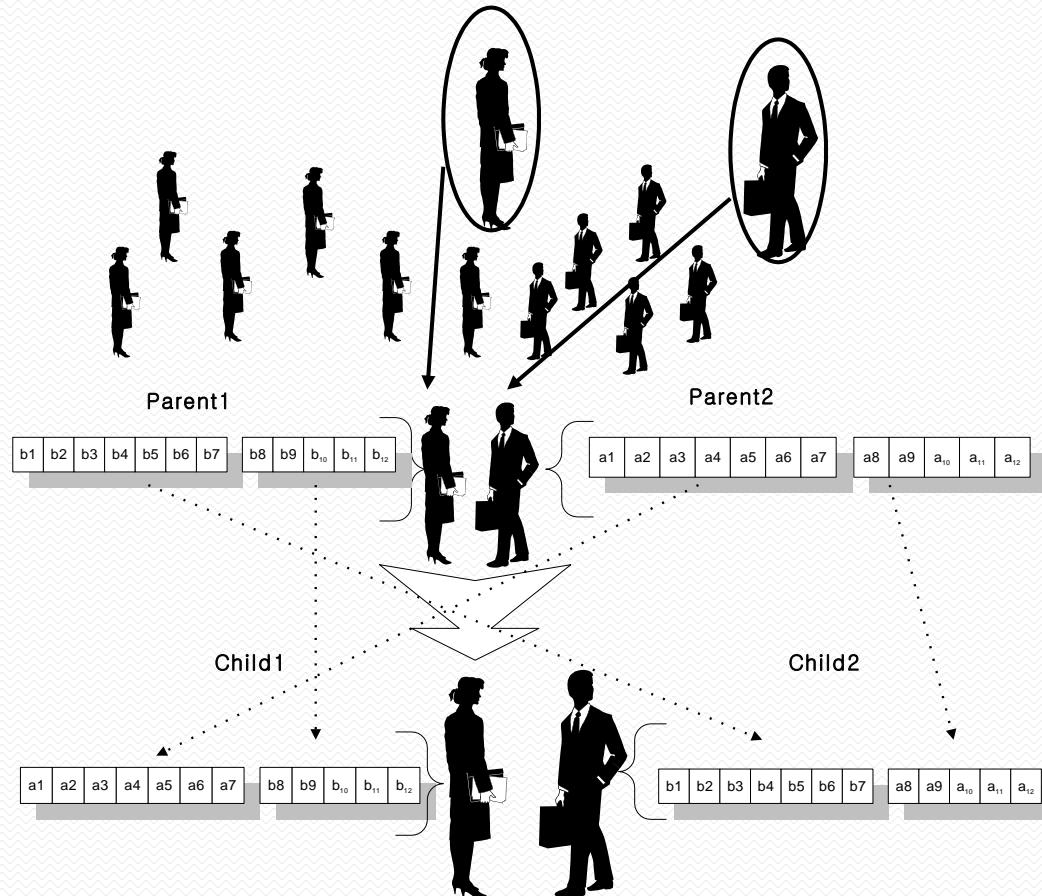
A book of "Adaptation"



1975

Genetic Algorithm : 존 홀랜드

- 저서 'Adaptation in Natural and Artificial Systems' 를 통해 제안
- 생물의 진화를 모방한 진화 연산의 대표적인 기법



HISTORY OF A.I.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^n P(B|A_j)P(A_j)}$$

베이지안 이론

- 베이지안 이론은 두 종류의 조건부 확률 사이의 관계를 정의함
- 다양한 분야에서 활용되는 대표적인 기법

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Likelihood **Class Prior Probability**
Posterior Probability **Predictor Prior Probability**

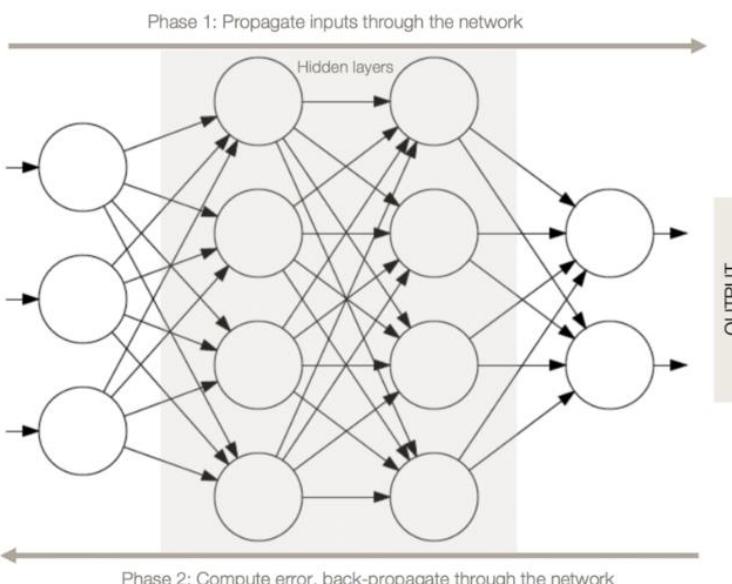
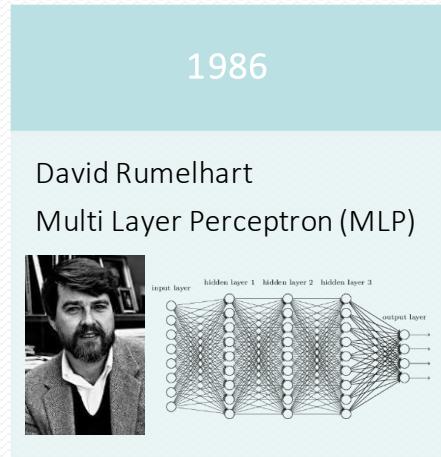
P(A), 사전확률(prior probability) : 결과가 나타나기 전에 결정되어 있는 A(원인)의 확률.

P(B|A), 우도확률(likelihood probability) : A(원인)가 발생하였다는 조건하에서 B(결과)가 발생할 확률.

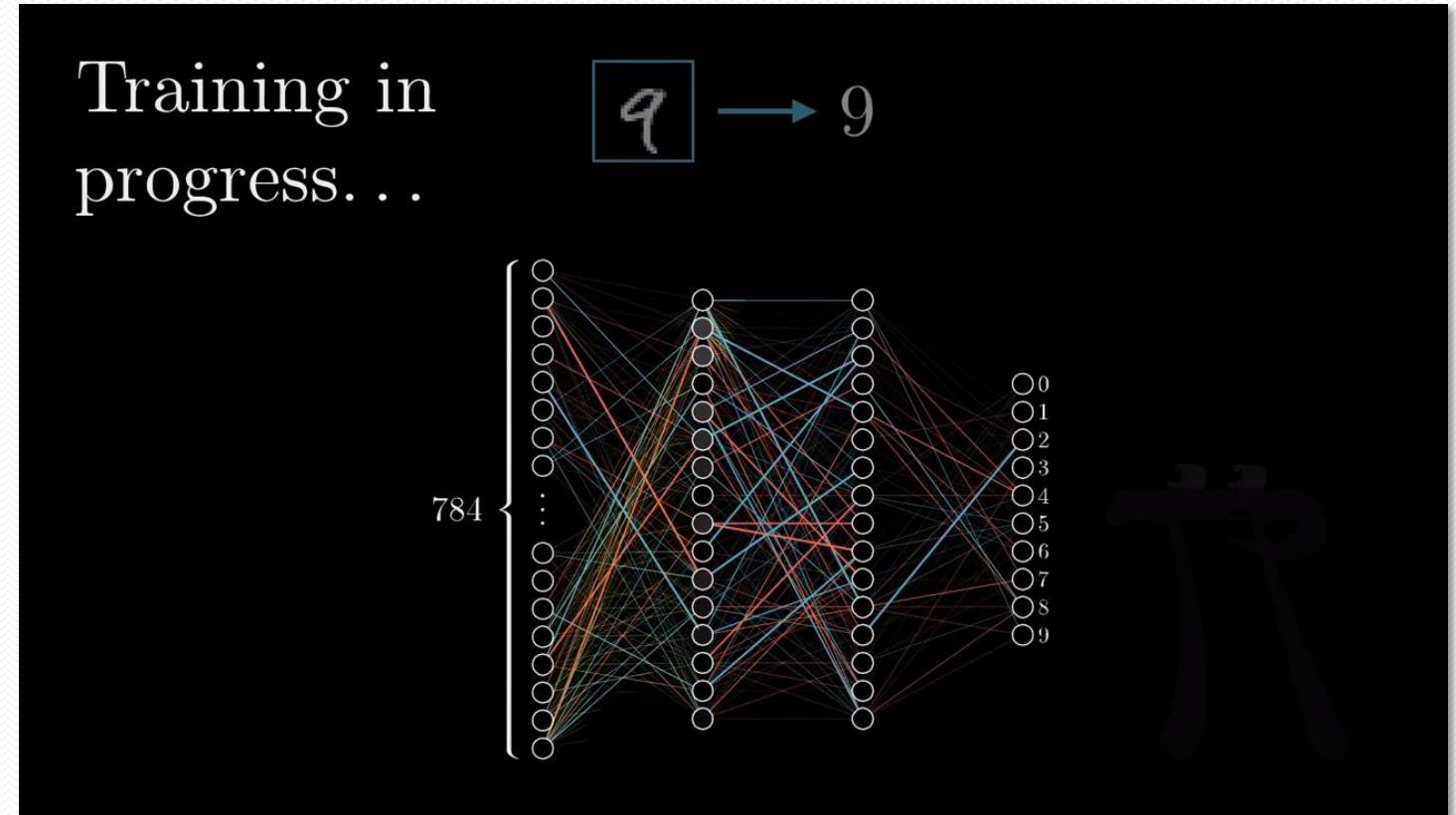
P(A|B), 사후확률(posterior probability) : B(결과)가 발생하였다는 조건하에서 A(원인)가 발생하였을 확률.

$$P(\text{질병}|\text{양성}) = \frac{P(\text{양성}|\text{질병})P(\text{질병})}{P(\text{양성})} = \frac{0.99 \times 0.01}{0.01 \times 0.99 + 0.99 \times 0.10} = 0.091$$

HISTORY OF A.I.



- Geoffrey Hinton, Ronald Williams 와 함께 **Back Propagation**에 대한 논문을 통해 Multi Layer Perceptron을 Training 할수 있음을 언급
- 다시 AI에 대한 열기를 가져옴



HISTORY OF A.I.

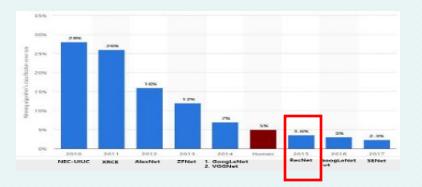
Yann Lecun
LeNet-5



INPUT 32x32
C1: 16 maps 8x8x28
S2: f maps 64@10x10
C3: f maps 16@10x10
S4: f maps 64@14x14
Convolutions Subsampling Conn

1998

ResNet
Accuracy is higher than human

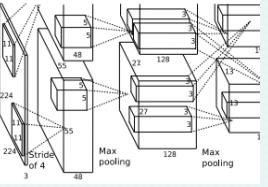


Model	Accuracy (%)
Net-5	25.1%
LeNet-5	26.0%
Baseline	19.0%
2011	17.0%
2012	15.4%
2013	13.1%
2014	11.3%
2015	10.5%
ResNet	15.7%
Baseline	10.5%
2016	10.5%

2015

2012

Alex Krizhevsky
AlexNet



2016

Google Brain
AlphaGo



2번째 위기

- Hidden layer가 깊어지면 back propagation이 제대로 안됨

LeNet-5

- Paper : Gradient-based learning applied to document recognition
- Multilayer Neural Network의 한계를 극복하기 위해 CNN을 활용
- MNIST Dataset

AlexNet

- Paper: ImageNet Classification with Deep Convolutional Neural Networks
- 2012년 ILSVRC 우승

ResNet

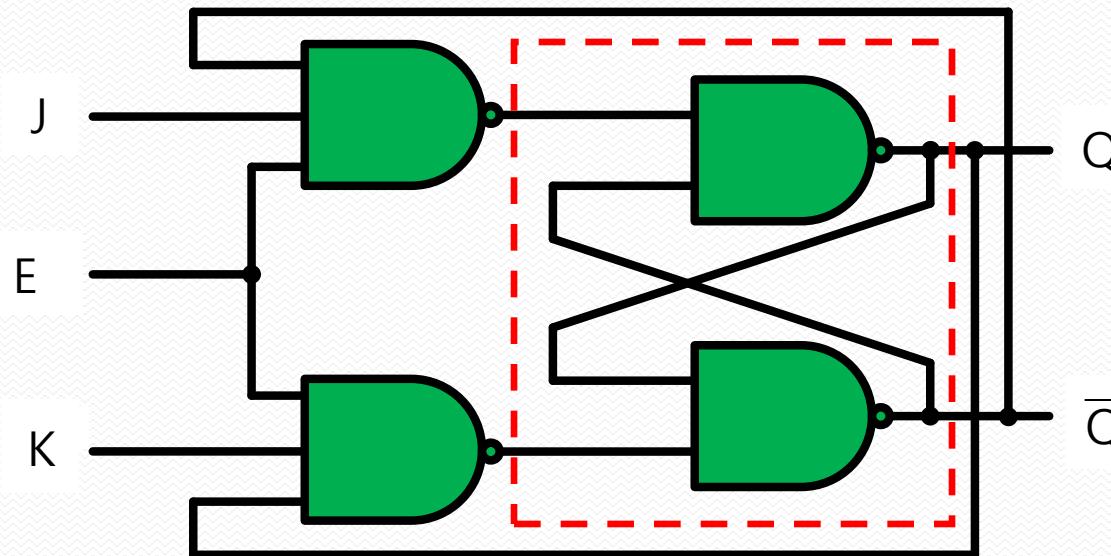
- 2015년 ILSVRC 우승
- Top-5 test error : 3.57%
- Residual neural network

HISTORY OF A.I. - 최근 언어 모델의 발전



Flip-flop (순차회로)

[NAND 게이트를 이용한 J-K 플립플롭]



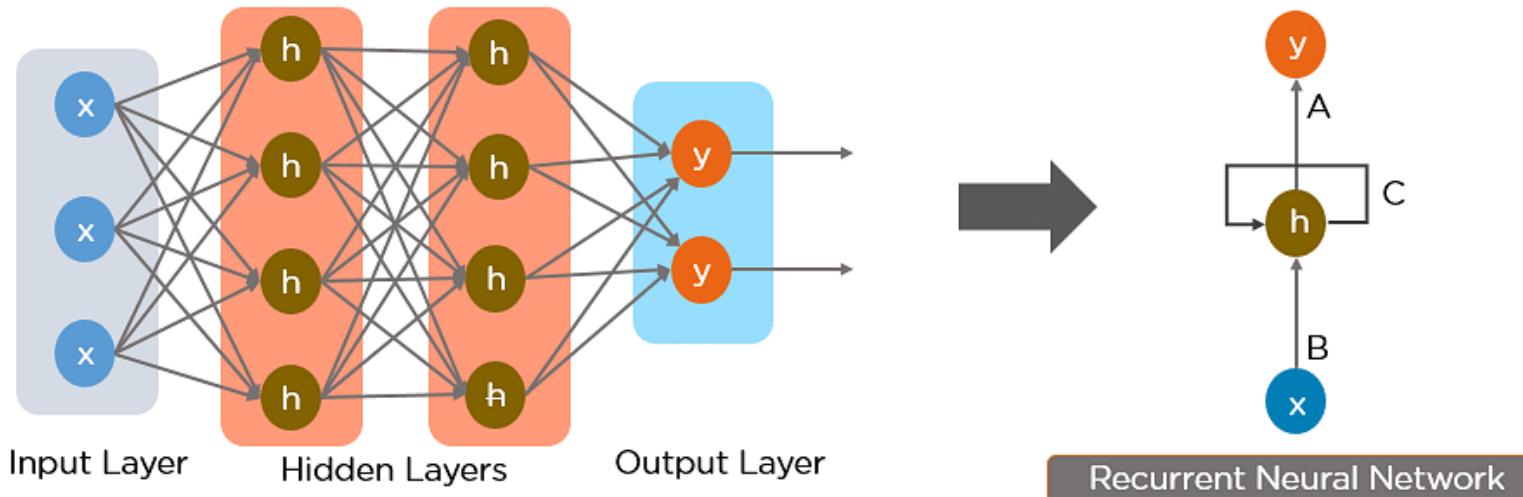
입력			출력
E	J	K	$Q(t+1)$
0	x	x	$Q(t)$ (불변)
1	0	0	$Q(t)$ (불변)
1	0	1	0
1	1	0	1
1	1	1	$Q(t)^*(toggle)$

- 순차회로를 이용해 기존의 조합회로에서 구현할 수 없었던 레지스터나 계수기를 구현 가능하게됨.
- RNN을 통해 기존 Neural Network에서 한계를 드러내었던 시계열 분석에 대한 성능향상이 진행됨.

순환 신경망(RNN)

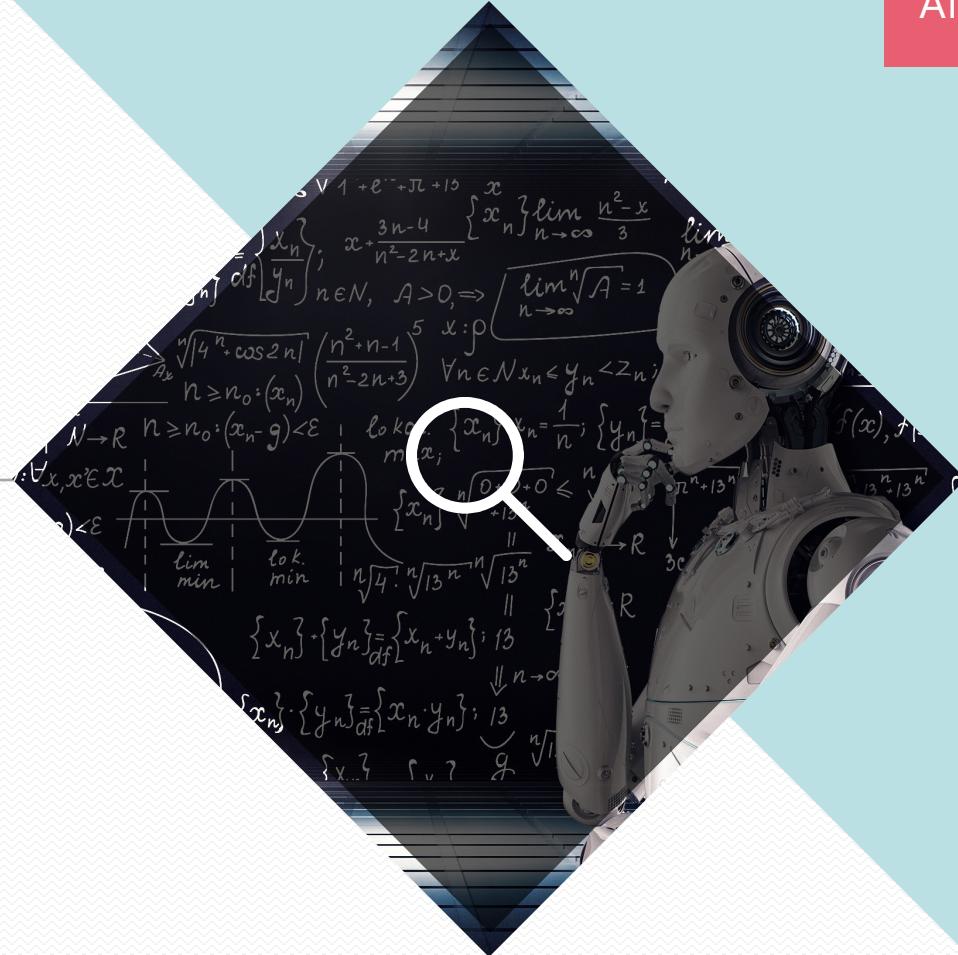
시간 개념을 반영할 수 있는 인공 신경망

- 시퀀스의 길이에 관계없이 입력과 출력을 받아들일 수 있는 네트워크 구조



ANN

ARTIFICIAL NEURAL NETWORKS



CONTENTS

NEURAL NETWORKS

- WHAT IS THIS?
- NEURON / WEIGHTS / BIASES

PERCEPTRON

- LOGICAL GATE EXAMPLE
- MULTI-LAYER PERCEPTRON , MLP

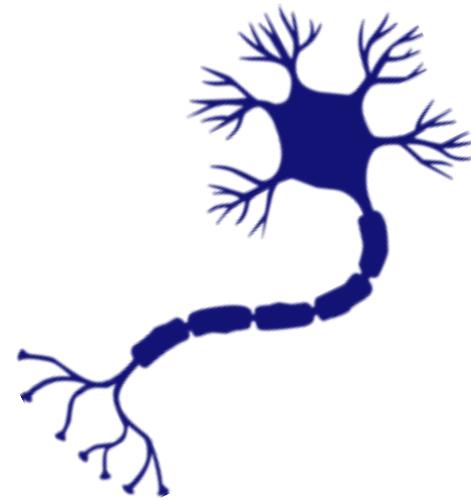
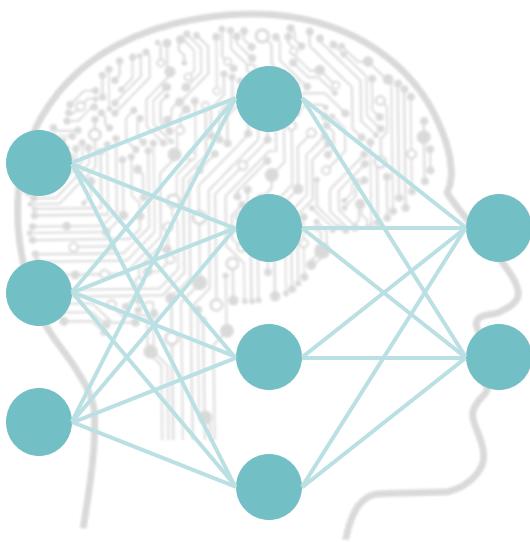
TECHNICAL DEEP DIVE

- ACTIVATION FUNCTIONS
- Optimizer

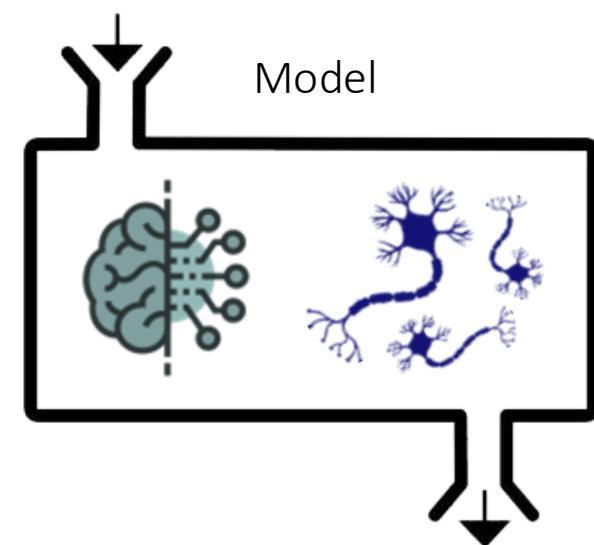
Hands-on

NEURAL NETWORKS OVERVIEW

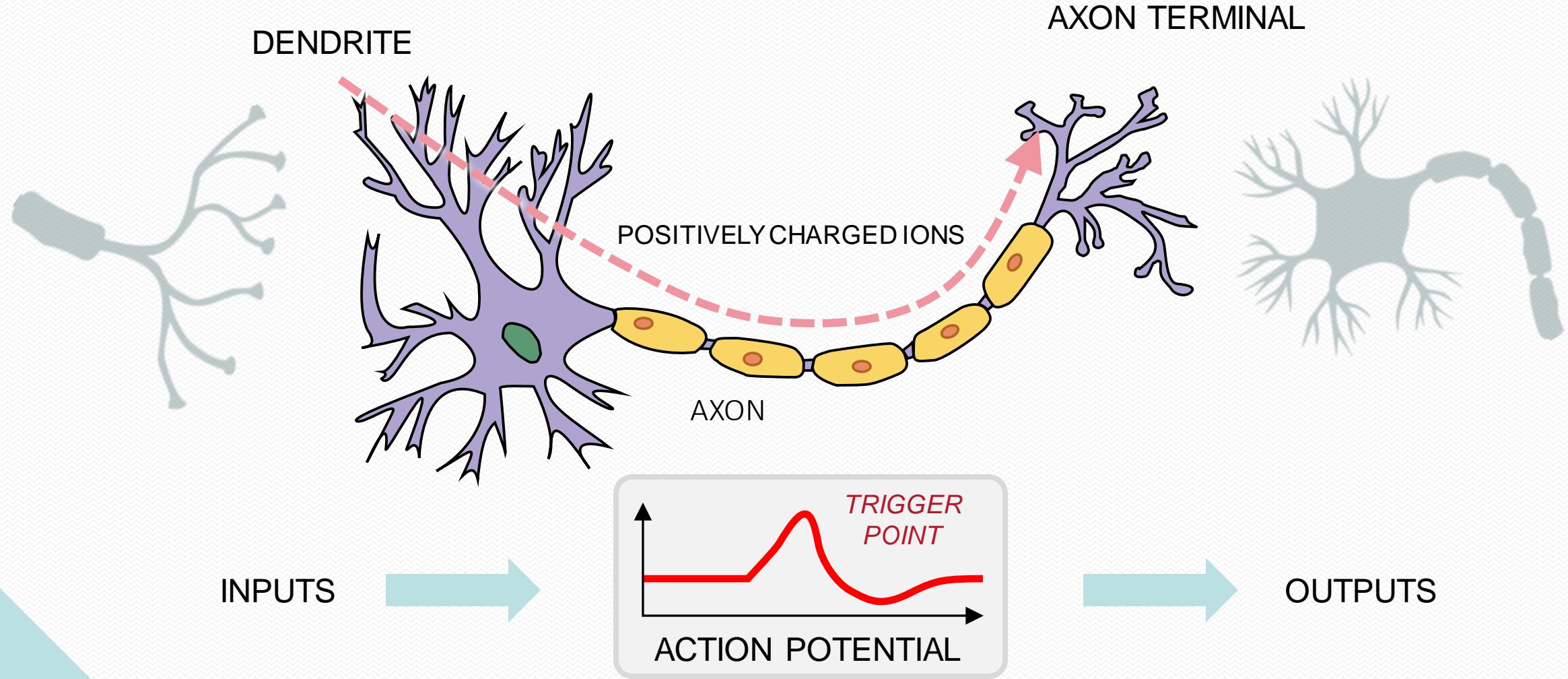
Neural networks



Neuron in human brain



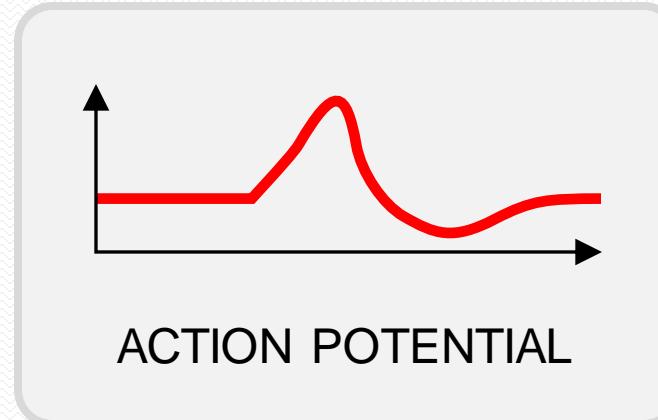
NEURON IN OUR BRAINS



ACTIVATION

NEURON

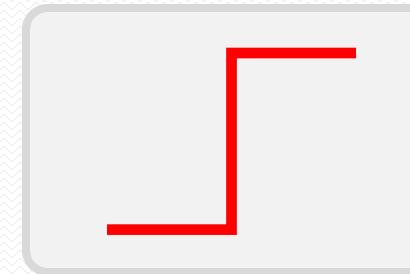
INPUTS



OUTPUTS

PERCEPTRON

INPUTS



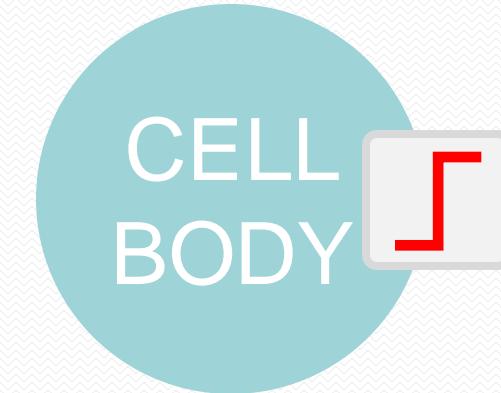
ACTIVATION FUNCTION

OUTPUT(S)

WEIGHTS AND BIASES

WEIGHTS

- CONTROL THE SIGNAL BETWEEN TWO NEURONS
- DETERMINES HOW THE INPUT AFFECTS THE OUTPUT



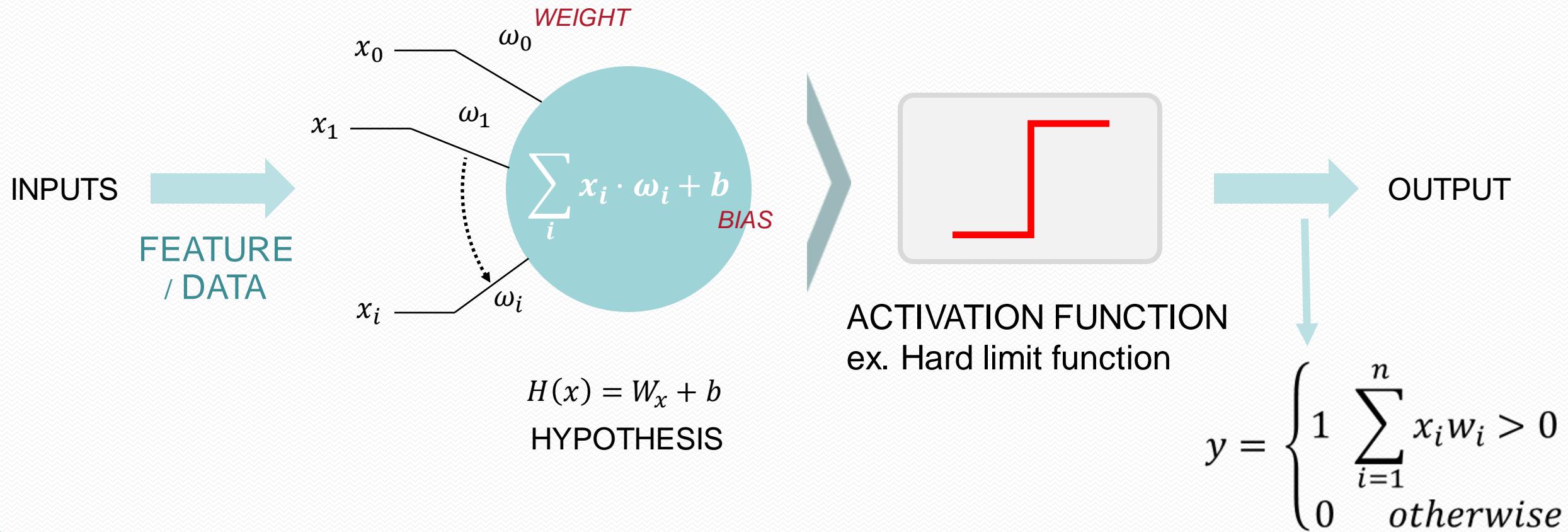
$$Y = \text{Activation}(\sum(\text{weight} \cdot \text{input}) + \text{bias})$$

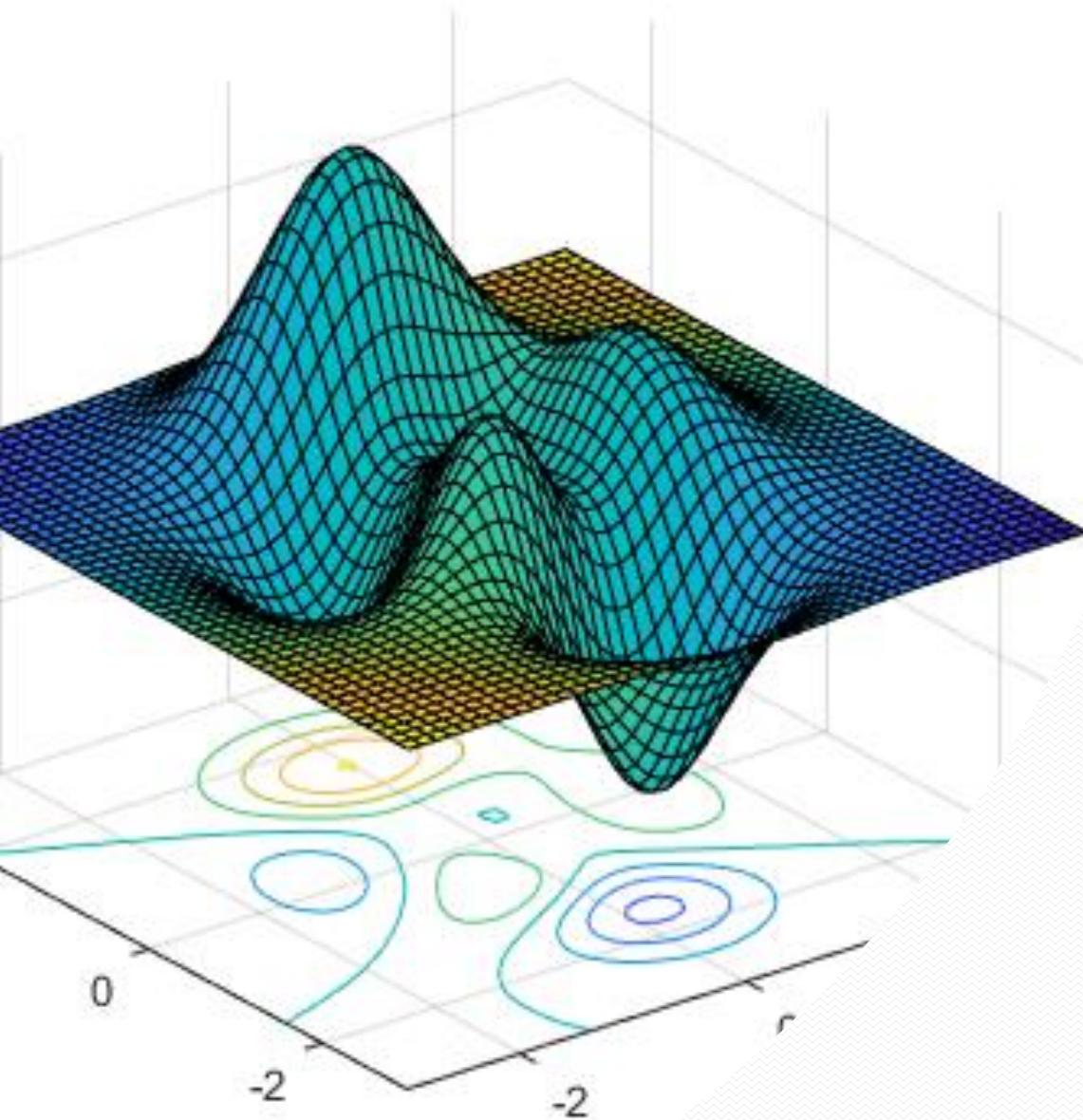
BIASES

- ADDITIONAL INPUT CONSTANT VALUES FOR NEXT NEURONS
- NOT INFLUENCED BY THE PREVIOUS NEURONS

These are the connections that go out of their own weights

PERCEPTRON

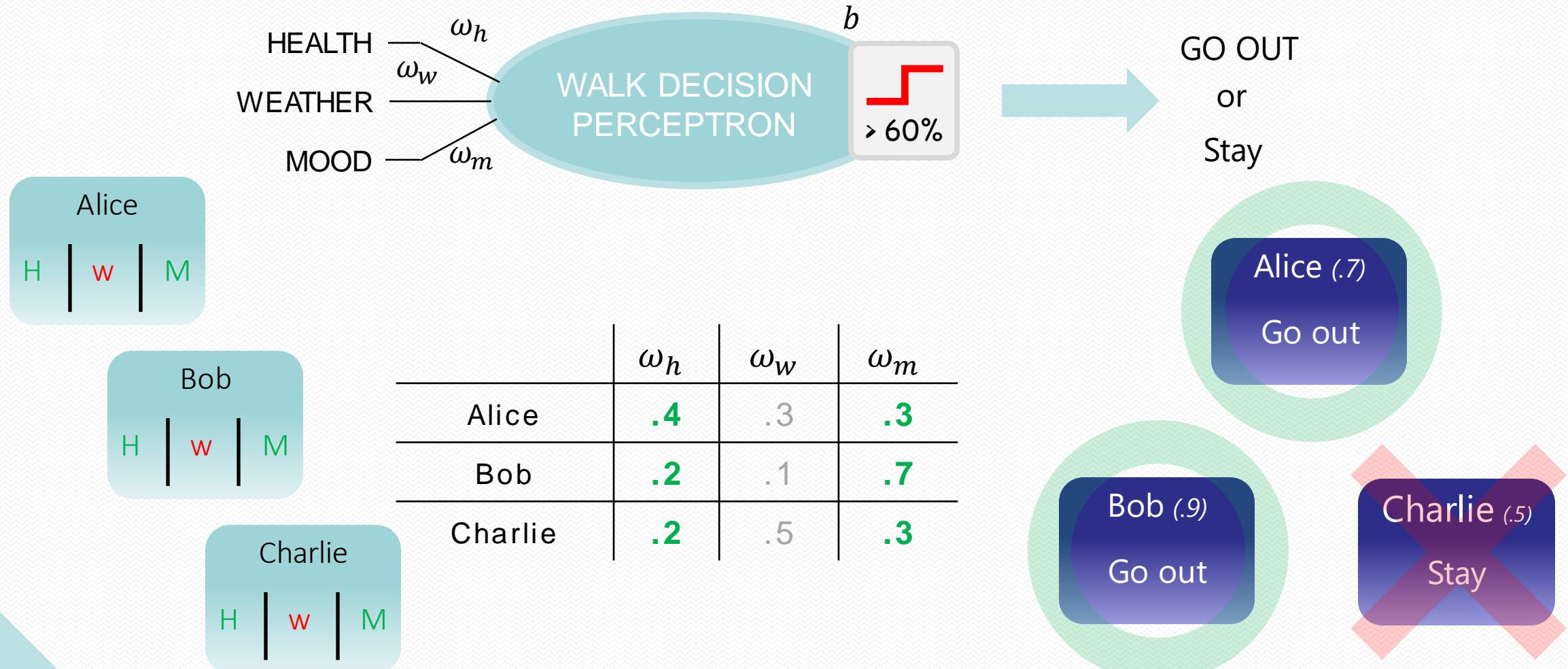




Perceptron

What a Perceptron can do?

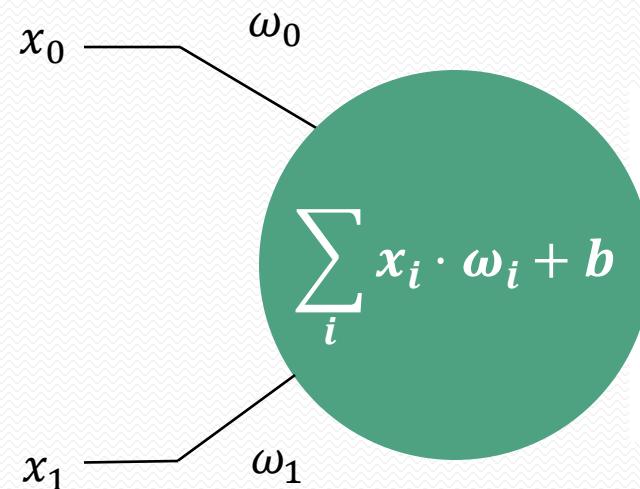
PERCEPTRON EXAMPLE



Perceptron Example – Logical OR

- Q. Find w_0 , w_1 and b that makes logical OR output.

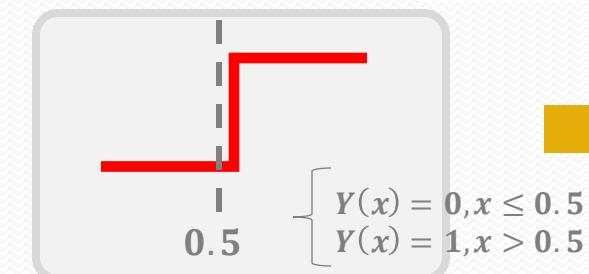
x_0	x_1	w_0	w_1	b	$H(x)$	Output
0	0	0.2	0.2	0.4	0.4	0
0	1				0.6	1
1	0				0.6	1
1	1				0.8	1



⋮

w_0	w_1	b
0.1	0.1	0.5
0.2	0.2	0.5
0.2	0.3	0.4
0.3	0.2	0.4
⋮	⋮	⋮

x_0	x_1	Output
0	0	0
0	1	1
1	0	1
1	1	1



Activation Function

output

Perceptron Example – Logical NAND

- Q. Find w_0 , w_1 and b that makes logical NAND output.

x_0	x_1	w_0	w_1	b	$H(x)$	Output
0	0	-0.5	-0.5	1.1	1.1	1
0	1				0.6	1
1	0				0.6	1
1	1				0.1	0

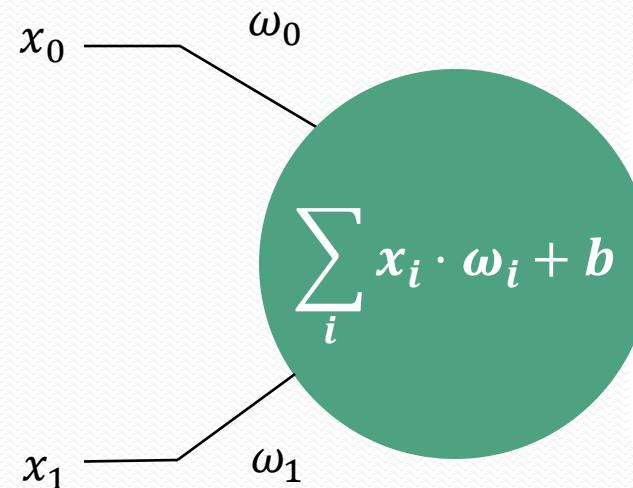
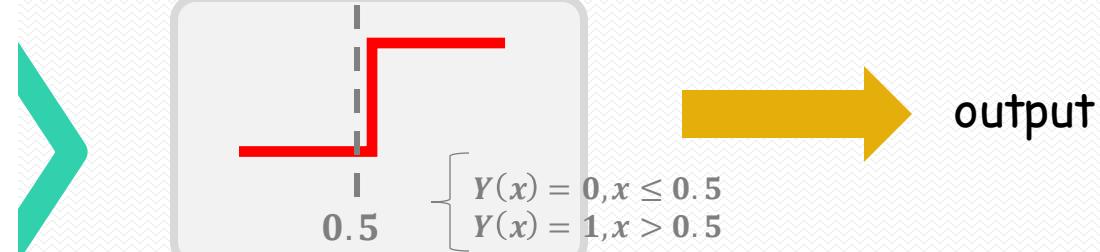


Table showing multiple sets of parameters w_0 , w_1 , and b that result in the same output for the given input patterns:

w_0	w_1	b
-0.2	-0.2	0.8
-0.3	-0.3	0.9
-0.4	-0.4	1.0
-0.5	-0.5	1.1
.	.	.

x_0	x_1	Output
0	0	1
0	1	1
1	0	1
1	1	0

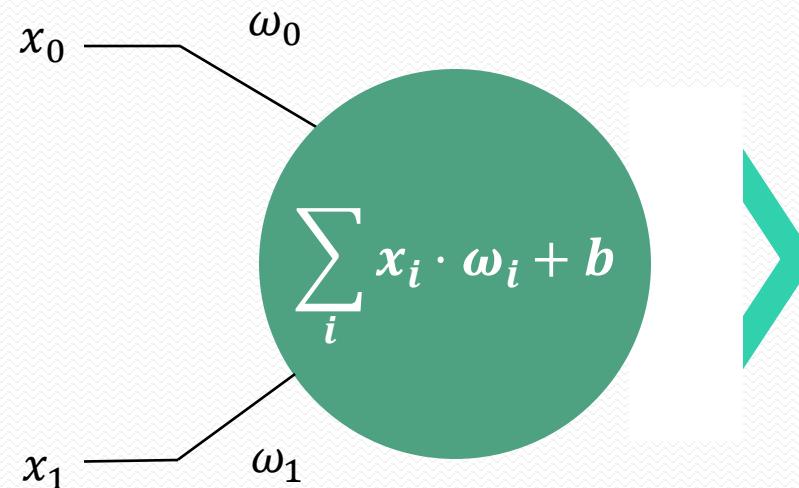


Activation Function

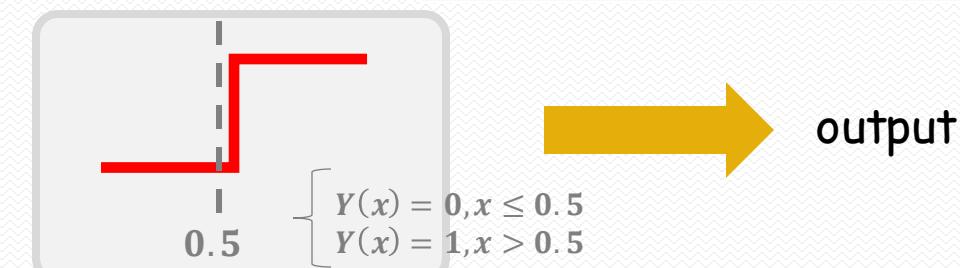
Perceptron Example – Logical XOR

- Q. Find w_0 , w_1 and b that makes logical XOR output.

x_0	x_1	w_0	w_1	b	$H(x)$	Output
0	0				0	0
0	1				0	0
1	0				0	0
1	1				0	0



x_0	x_1	Output
0	0	0
0	1	1
1	0	1
1	1	0



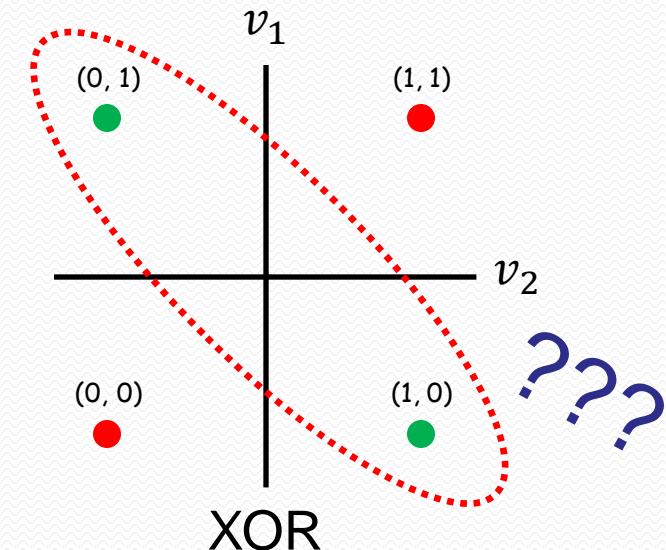
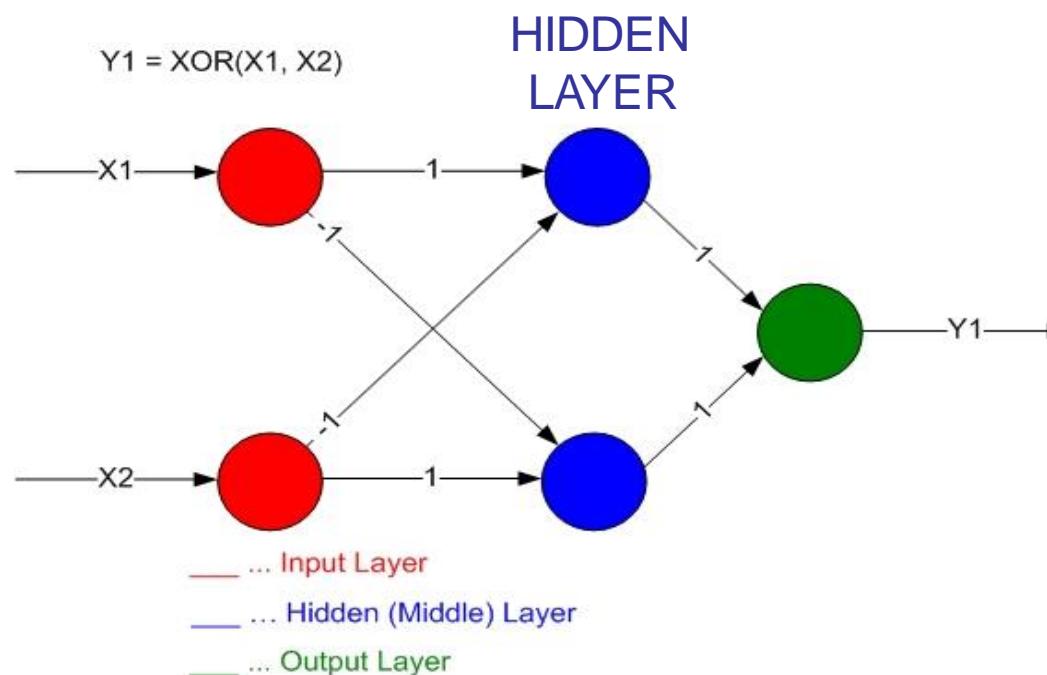
Activation Function

PROBLEM OF PERCEPTRON

PERCEPTRON IS NOT AN ALL-ROUND PLAYER

PROBLEMS THAT ARE NOT CLASSIFIED AS LINEAR CANNOT BE CALCULATED

e.g., **XOR** logic / impossible to solve with a single layer of perceptron

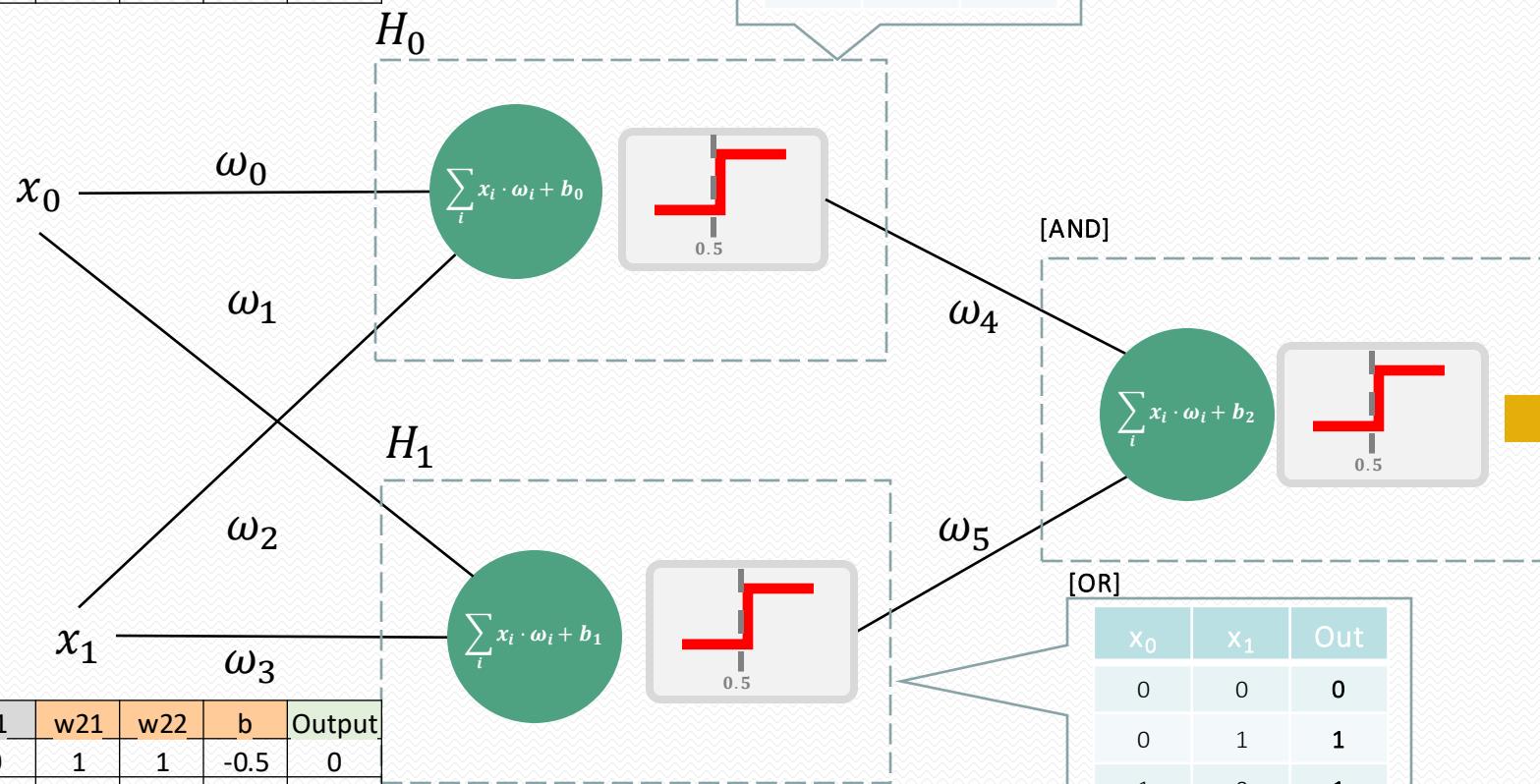


Multi Layer Perceptron - XOR

x_0	x_1	w_{11}	w_{12}	b	Output
0	0	1	1	-1.5	0
0	1				0
1	0				0
1	1				1

[NAND]

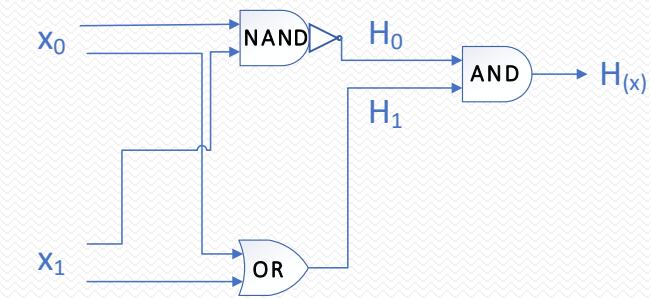
x_0	x_1	Out
0	0	1
0	1	1
1	0	1
1	1	0



x_0	x_1	w_{21}	w_{22}	b	Output
0	0	1	1	-0.5	0
0	1				1
1	0				1
1	1				1

[OR]

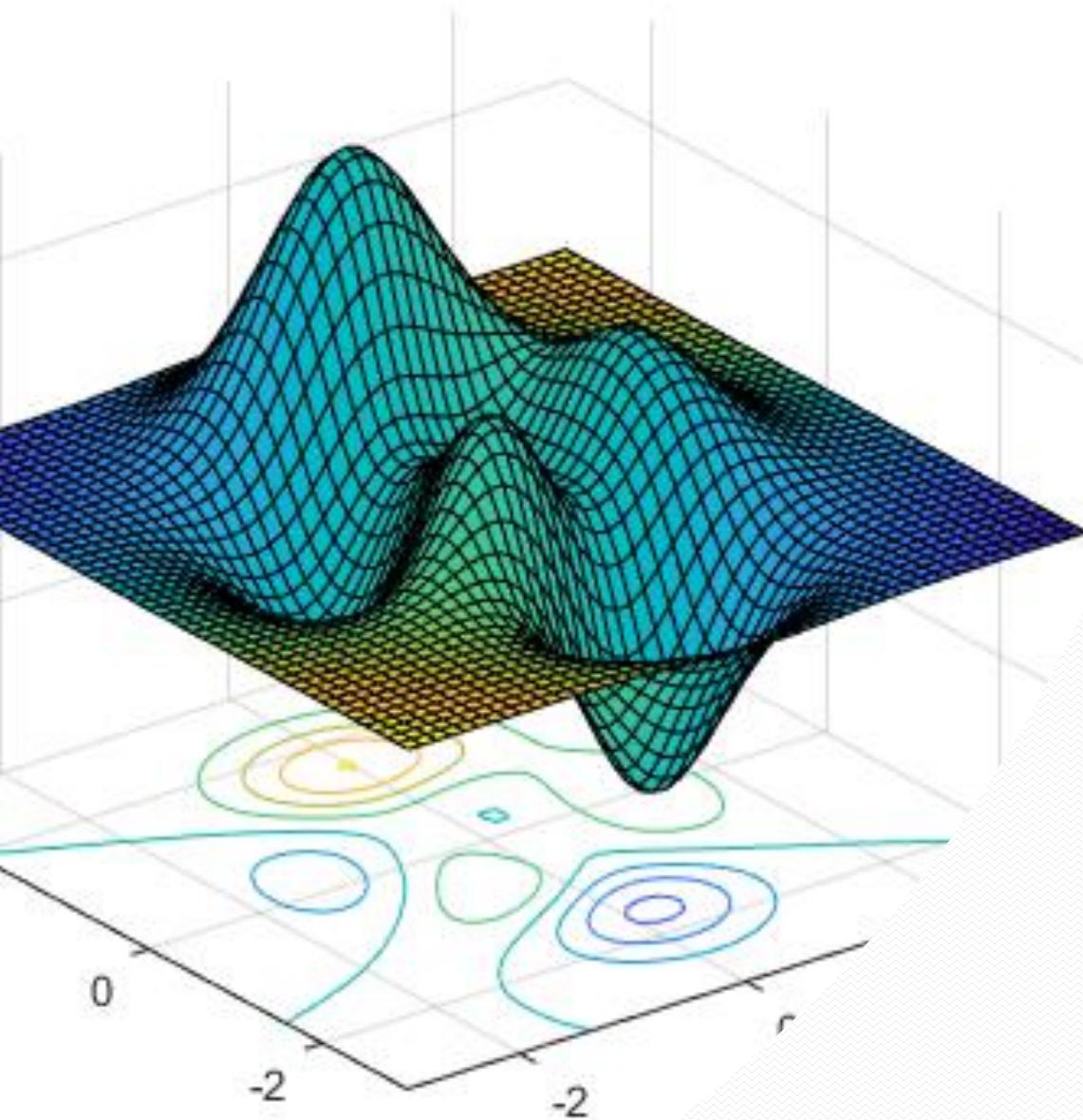
x_0	x_1	Out
0	0	0
0	1	1
1	0	1
1	1	1



[XOR]

w_0	x_1	Out
0	0	0
0	1	1
1	0	1
1	1	0

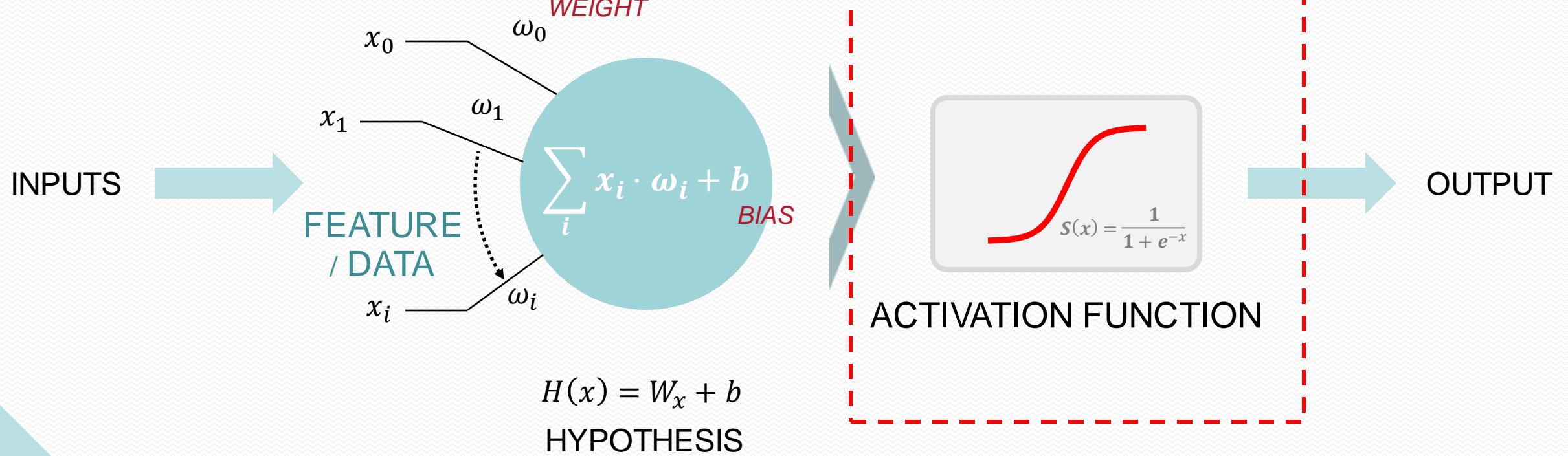
x_0	x_1	w_{11}	w_{12}	b	Output
0	0	-1	1	-0.5	0
0	1				1
1	0				1
1	1				0



Deep-dive

Activation Function

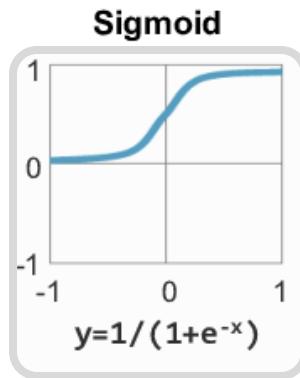
PERCEPTRON



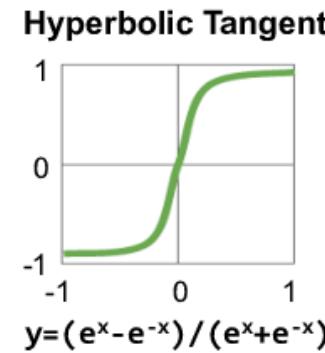
ACTIVATION FUNCTIONS

DEFINES THE OUTPUT OF A NODE USING A GIVEN INPUT OR SET OF INPUTS

**Traditional
Non-Linear
Activation
Functions**

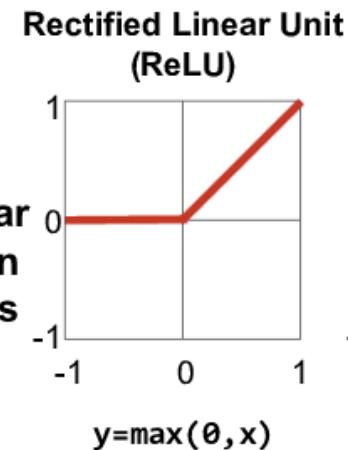


For binary classification

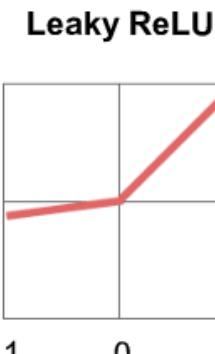


Widely used in hidden layers

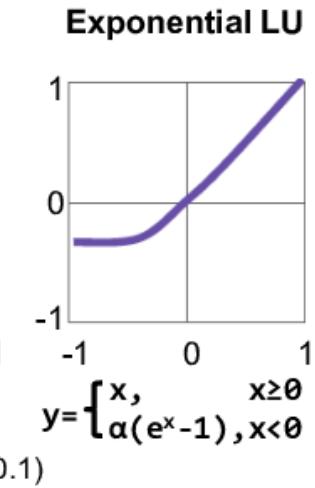
**Modern
Non-Linear
Activation
Functions**



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



$$y = \max(\alpha x, x)$$



$$y = \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$$

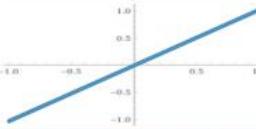
$$\alpha = \text{small const. (e.g. 0.1)}$$

ACTIVATION FUNCTIONS

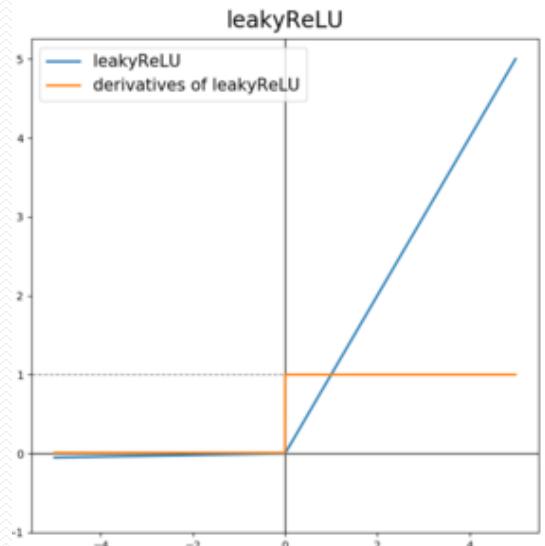
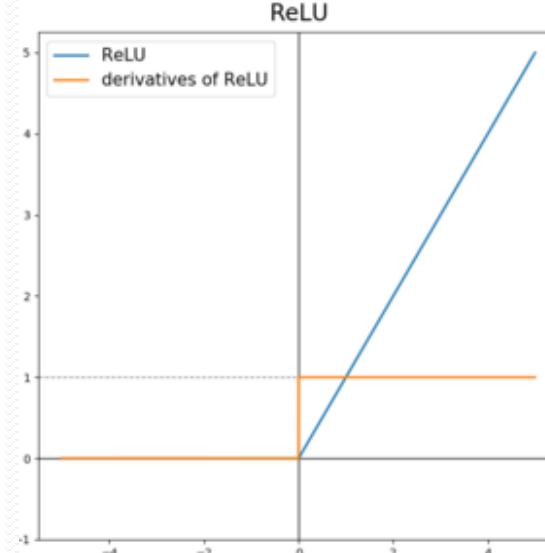
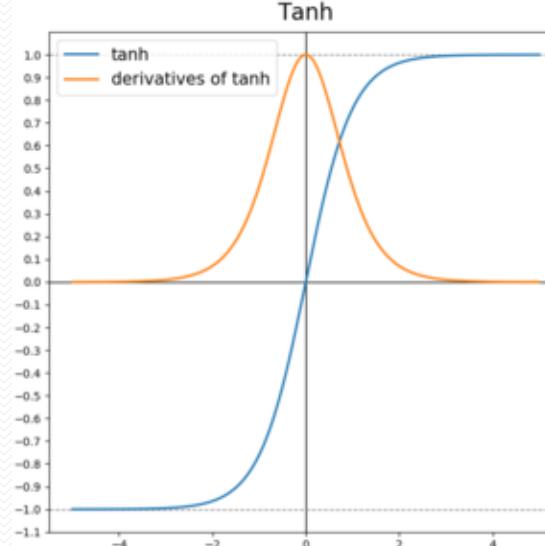
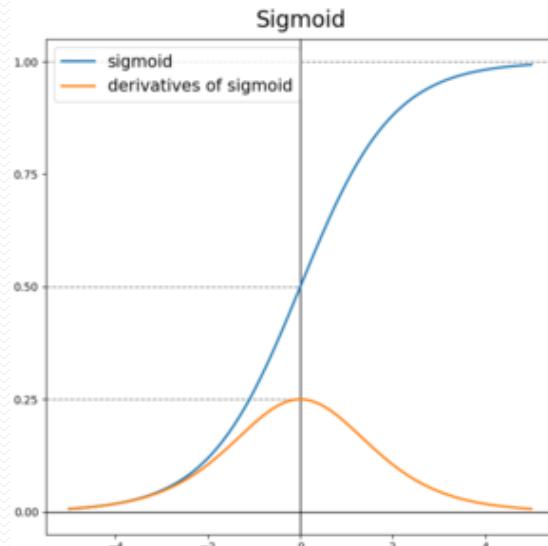
ACTIVATION FUNCTIONS

Why derivative/differentiation is used?

To know in which direction and how much to change or update the curve depending upon the slope.

Name	Plot	Equation	Derivative
Linear		$f(x) = x$	$f'(x) = 1$
Sigmoid		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Leaky ReLU		$f(x) = \begin{cases} ax & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} a & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

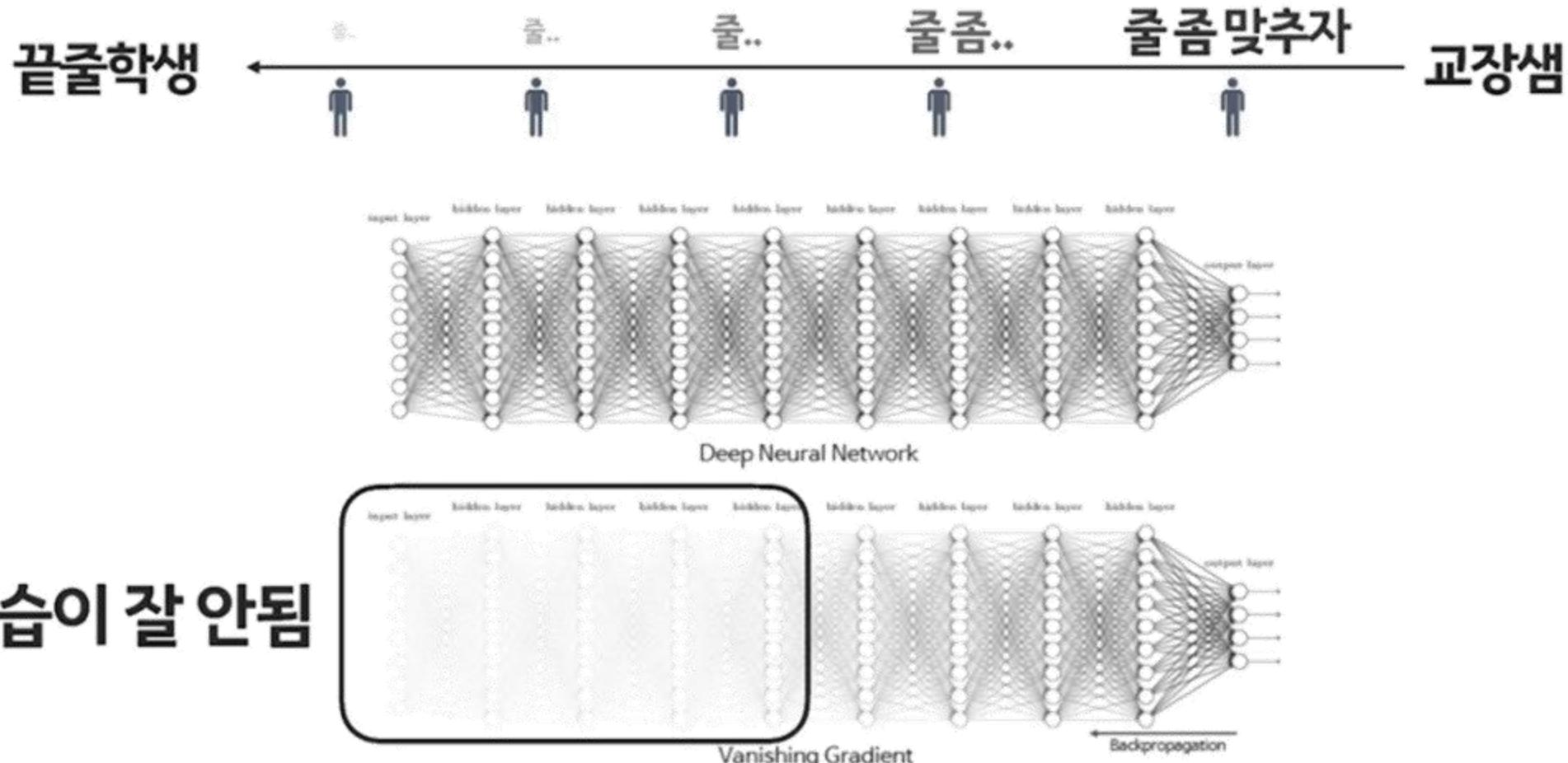
ACTIVATION FUNCTIONS



	Sigmoid	Tanh	ReLU	Leaky ReLU
Disadvantage	Vanishing Gradient Computational penalty	Vanishing Gradient	Dying ReLU problem	Result not consistent for negative
Advantage	1 or 0, binary value	Zero centered	Computational efficient Non-linear	Prevent negative Values become zero

Vanishing gradient

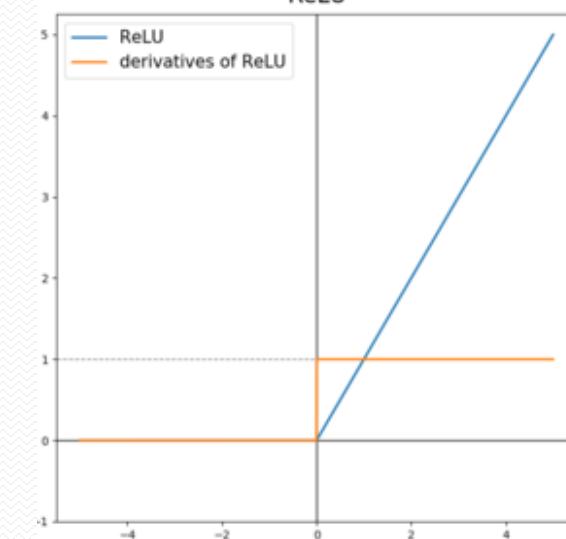
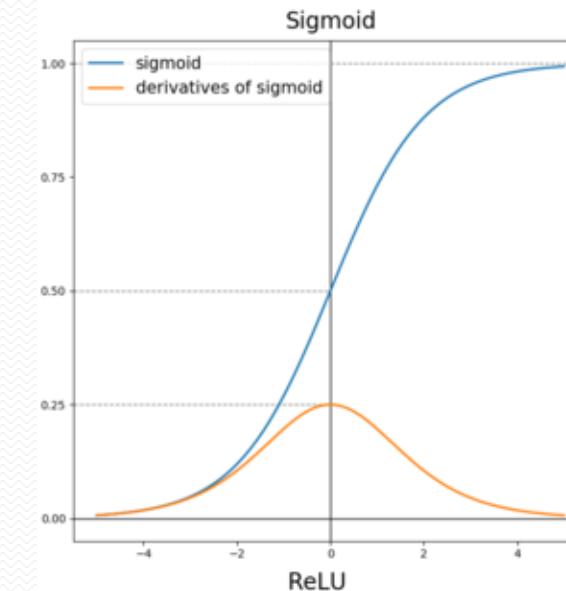
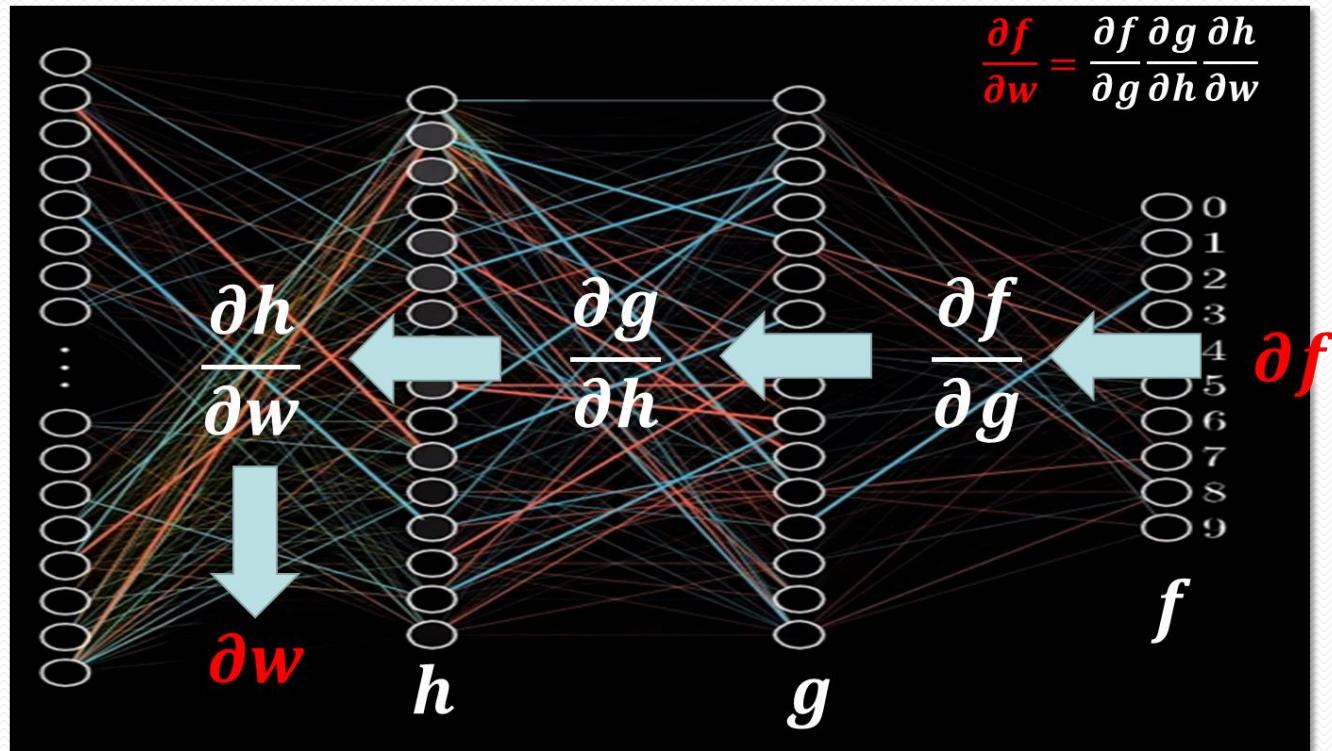
레이어가 깊을 수록 업데이트가 사라져감.



학습이 잘 안됨

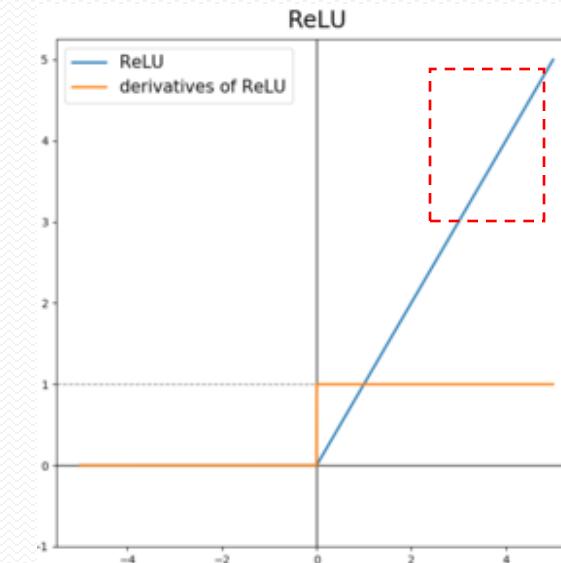
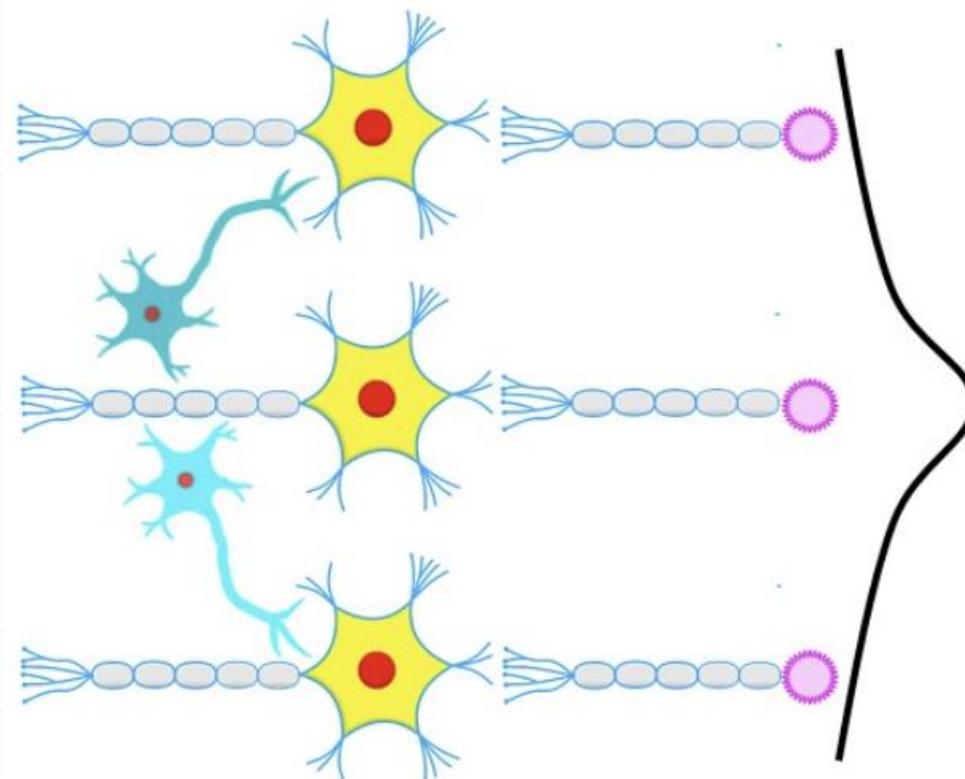
ACTIVATION FUNCTIONS: Vanishing Gradient

CALCULATE THE ∂w from ∂f

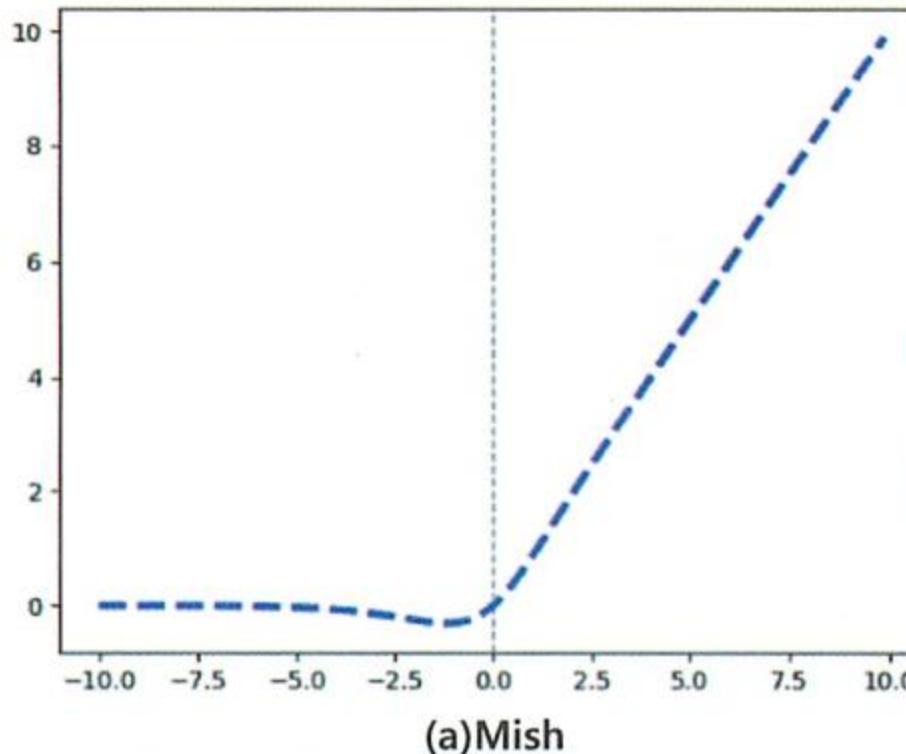


ACTIVATION FUNCTIONS: ReLu vs 측면 억제

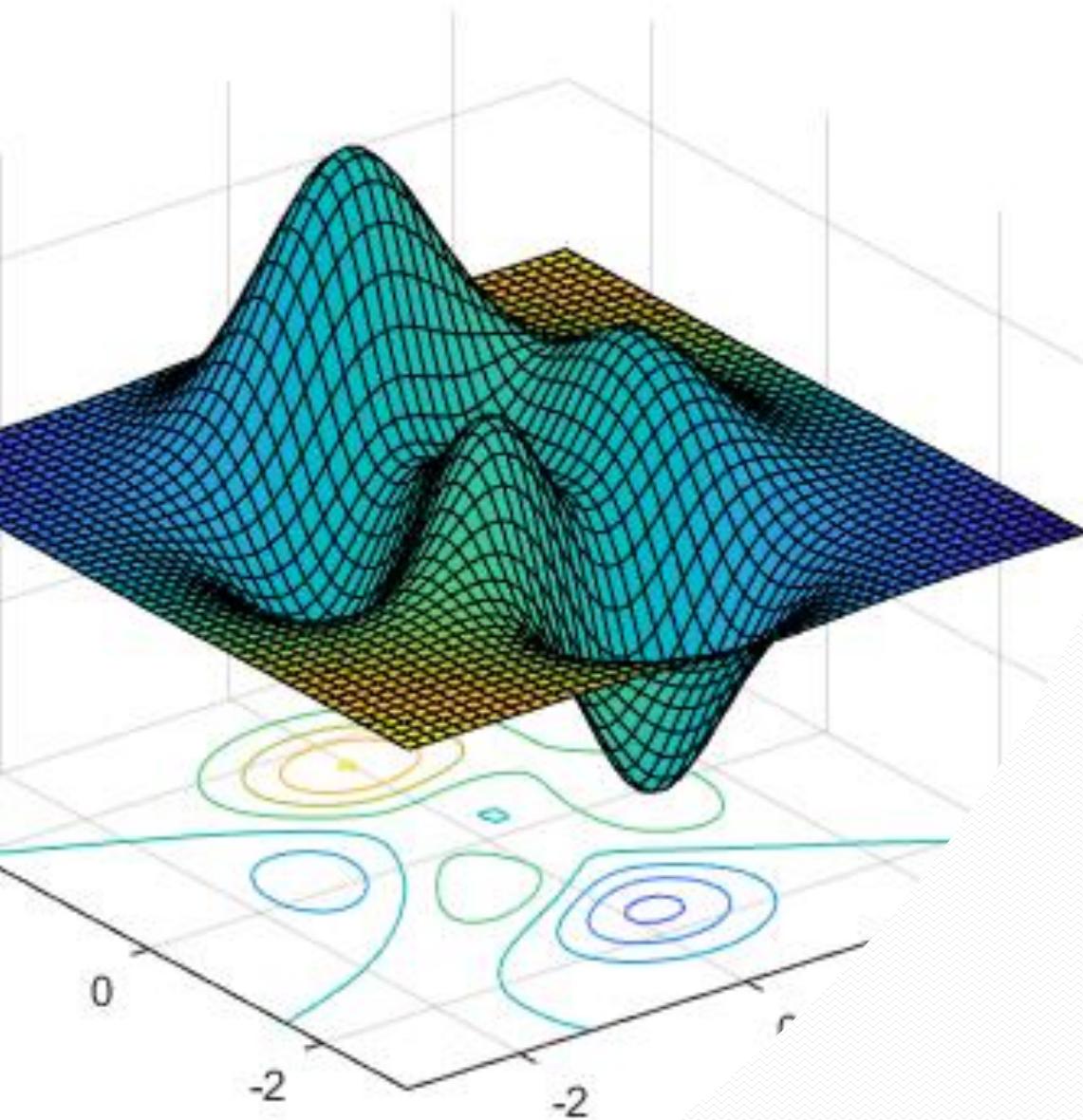
입출력단을 갖는 2개의 신경세포에서 각각의 신경세포의 입력 또는 출력이 서로 다른 신경세포의 입력 또는 출력에 의해 억제되는 현상.



ACTIVATION FUNCTIONS: Yolov4 Mish func.



```
def MISH(x):  
    return x * np.tanh(np.log(1 + np.exp(x)))
```

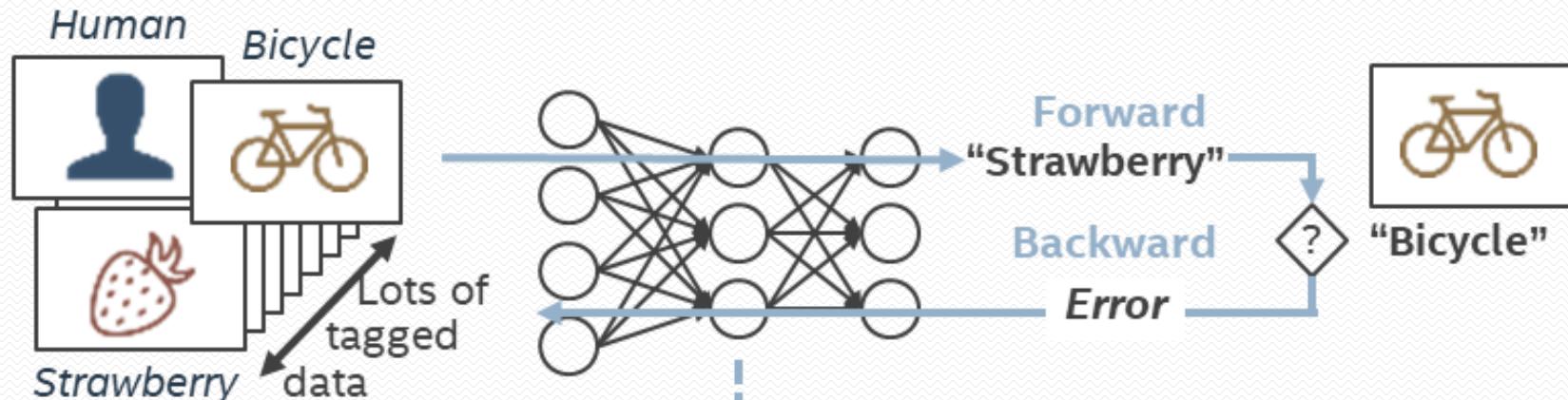


Deep-dive

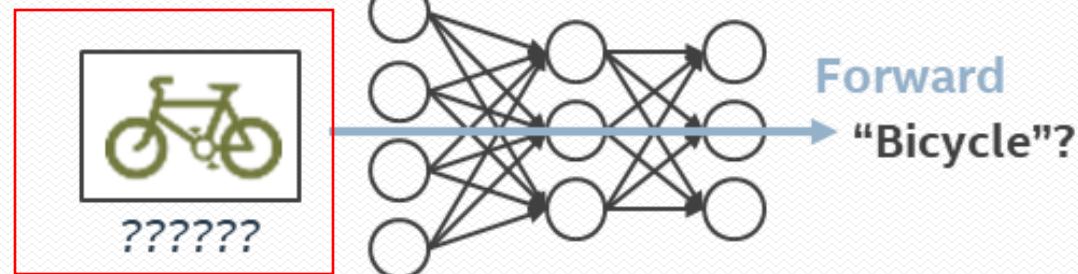
ANN: technical in details

TRAINING & INFERENCE

TRAINING:



INFERENCE:

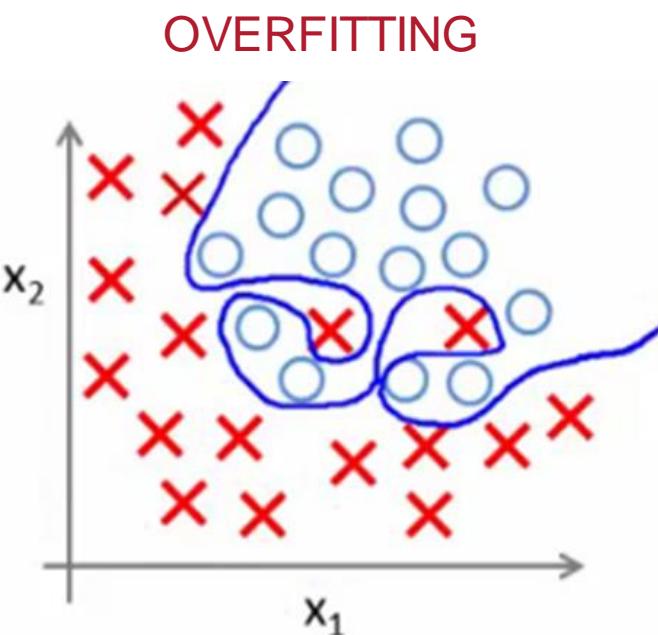
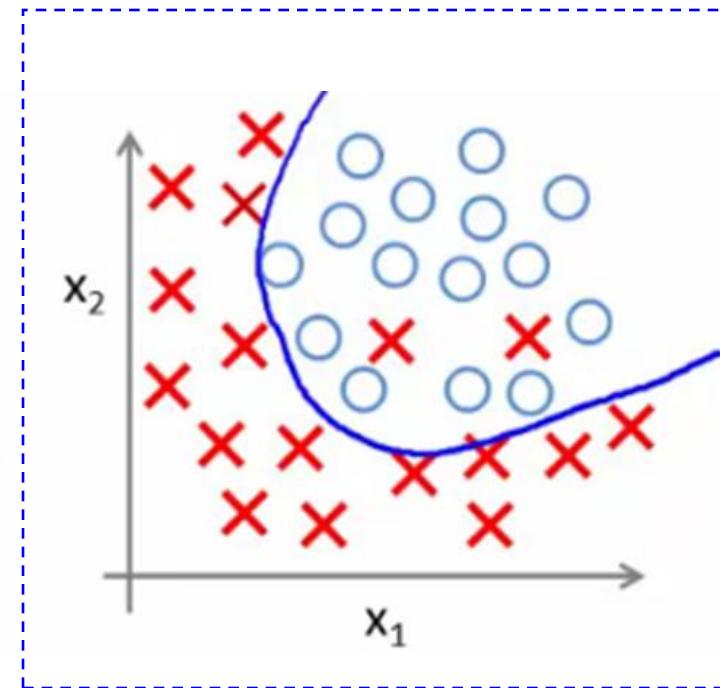
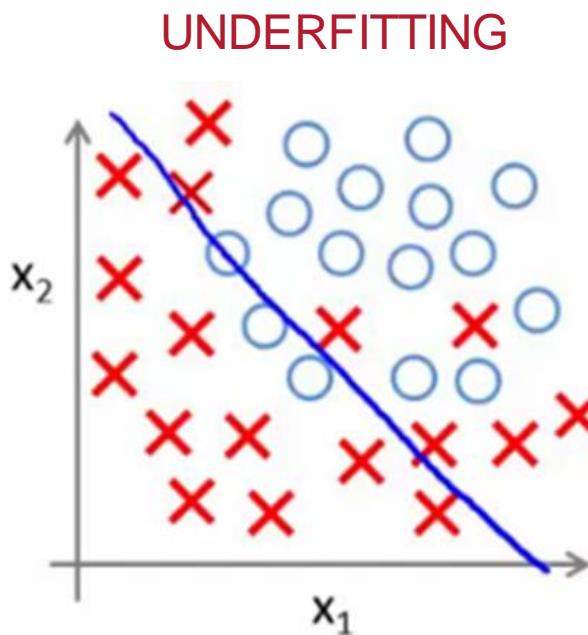


OVERFITTING / UNDERFITTING

HIGH COMPLEXITY MODEL = HIGH POSSIBILITY OF OVERFITTING

If the model is too specific, generalization is reduced

When new test data sets come in → they may not be well categorized



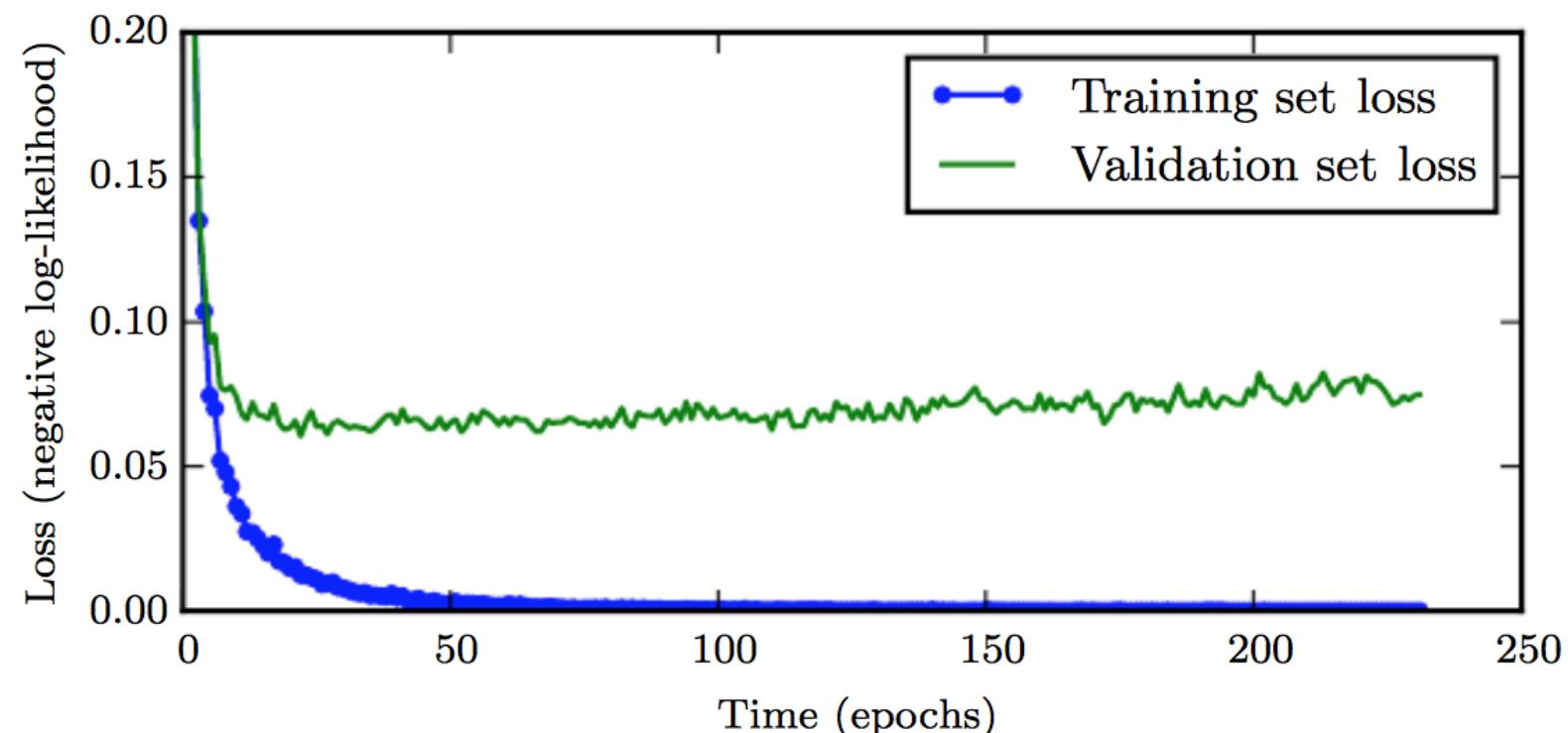
GENERALIZATION / RIZATION

EARLY STOPPING

Stop training when the gap between validation error and training error becomes large

Simple to do without additional hyperparameters, but mixing two problems

→ *UNDERFITTING AND OVERFITTING*

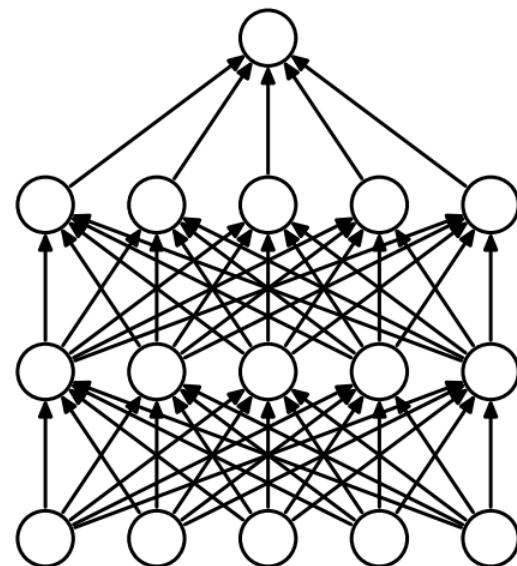


DROPOUT

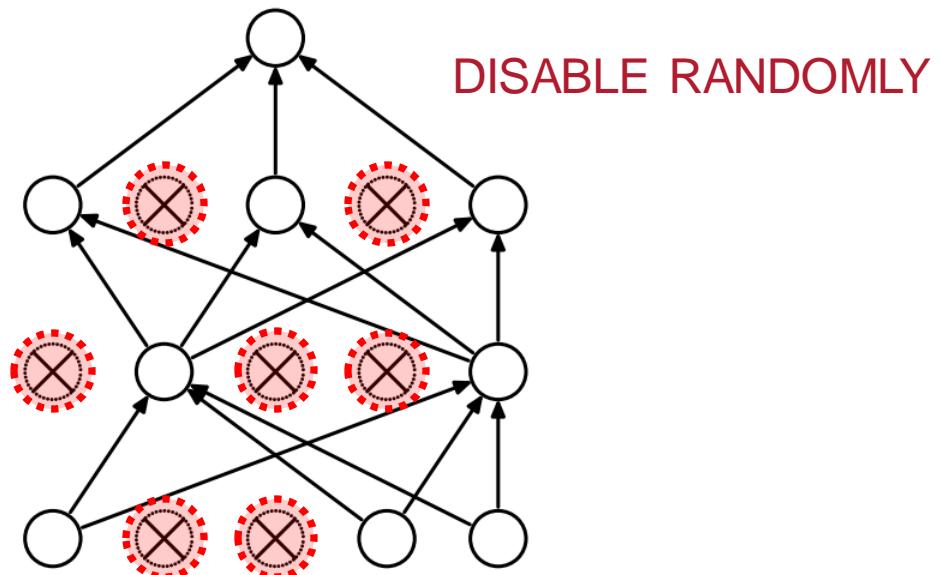
RANDOMLY SET SOME NEURONS TO ZERO IN THE FORWARD PASS

Voting effect → makes the result not biased

Anti co-adaptation → creates a robust network

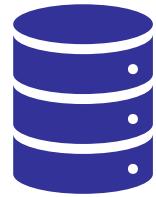


(a) Standard Neural Net



(b) After applying dropout.

DATASET



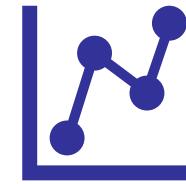
TRAIN

The actual dataset that we use to train the model (**weights** and **biases** in the case of a Neural Network).



VALIDATION

The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model **hyperparameters**.



TEST

The sample of data used to provide an **unbiased** evaluation of a final model fit on the training dataset.

HYPERPARAMETERS

OF HIDDEN LAYER

determines how many hidden layers is in neural network

OF EPOCH

finds the boundary value of the epoch where the error is minimized

BATCH SIZE

means the size of the data to be divided into when learning

LEARNING RATE

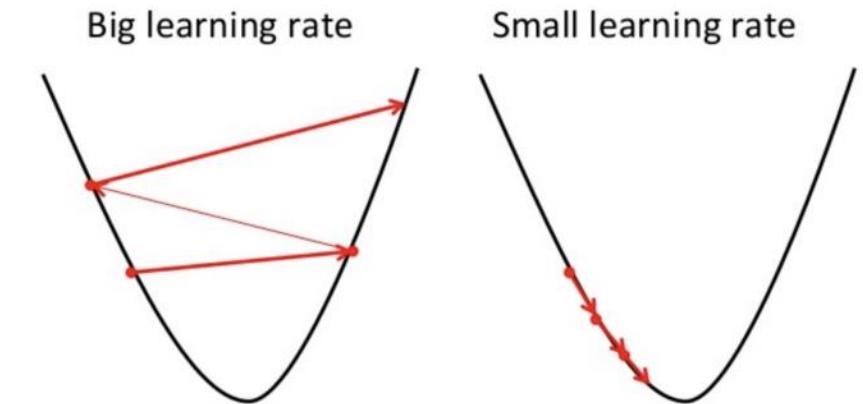
determines how fast or slow to move per iteration

INITIAL VALUE OF EACH PARAMETER

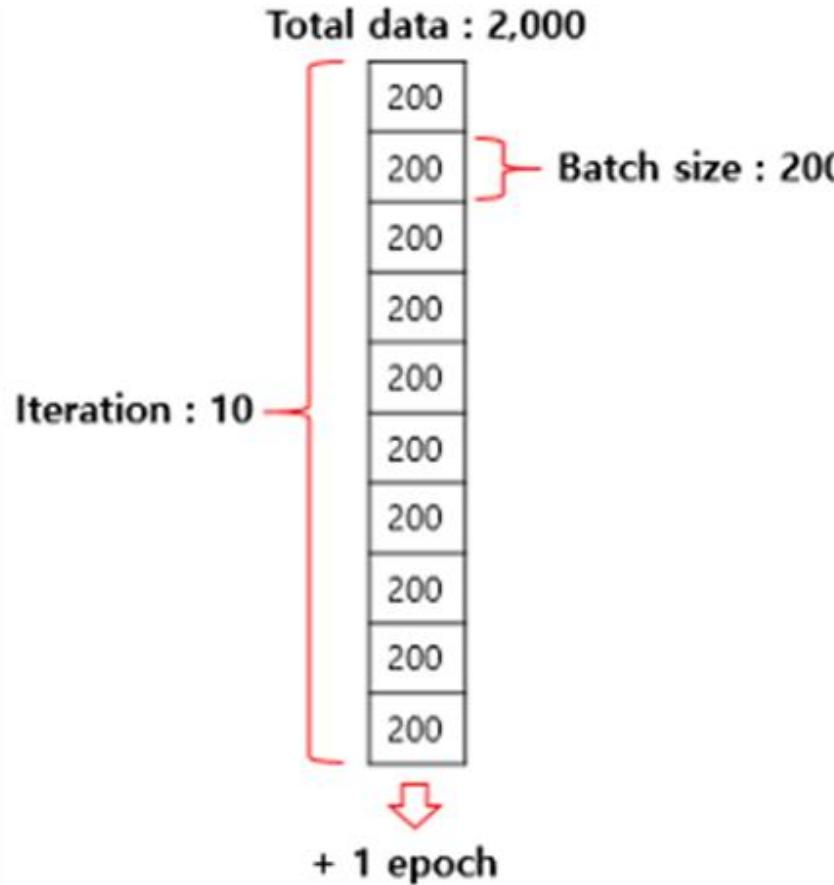
determines the initial value when running iteration

LOSS / COST FUNCTION

Decides how to calculate between actual result and calculated output



EPOCH & ITERATION & BATCH



- 데이터가 50,000개 있다.
- 1차원 선형회귀 문제 : $50,000 \times 4\text{bytes} = 200,000\text{bytes} = 0.2\text{MB}$
- MNIST 이미지 : $28 \times 28 \times 50,000 \times 4\text{bytes} = 156\text{MB}$
- 저화질 이미지 : $280 \times 280 \times 50,000 \times 4\text{bytes} = 15,600\text{MB} = 15.6\text{GB}$
- 중간화질 이미지 : $560 \times 560 \times 50,000 \times 4\text{bytes} = 15,600\text{MB} = 62.4\text{GB}$

EPOCH & ITERATION & BATCH

일부분의 데이터(batch)를 사용하여 Gradient를 추정하면 되지 않을까?

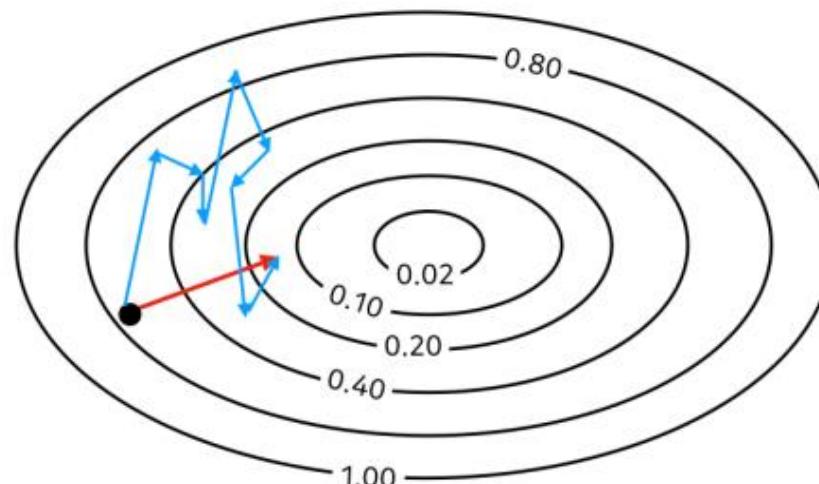
Total data : 2,000



If Batch size is Big(BGD)?



If Batch size is small(SGD)?



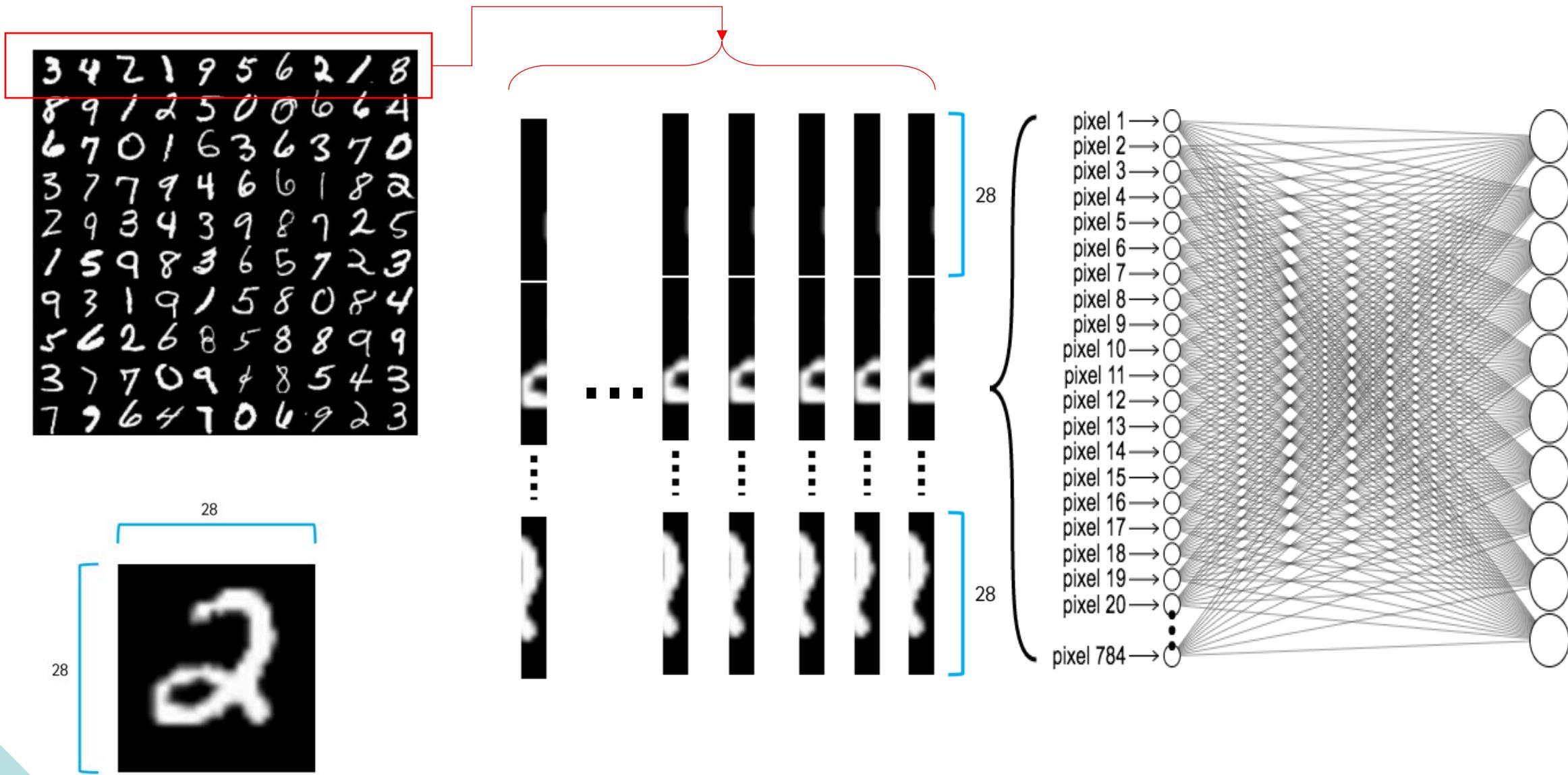
Then, Big batch size is more efficient?

→ requires too much memory

→ Too many parameters makes slow learning speed

- 계산속도 증가
- 메모리 용량 감소
- Local minimu을 빠져나갈 가능성 증가
- 부정확한 Gradient
- 가장 빠른 방향으로 Search방향 설정 불가

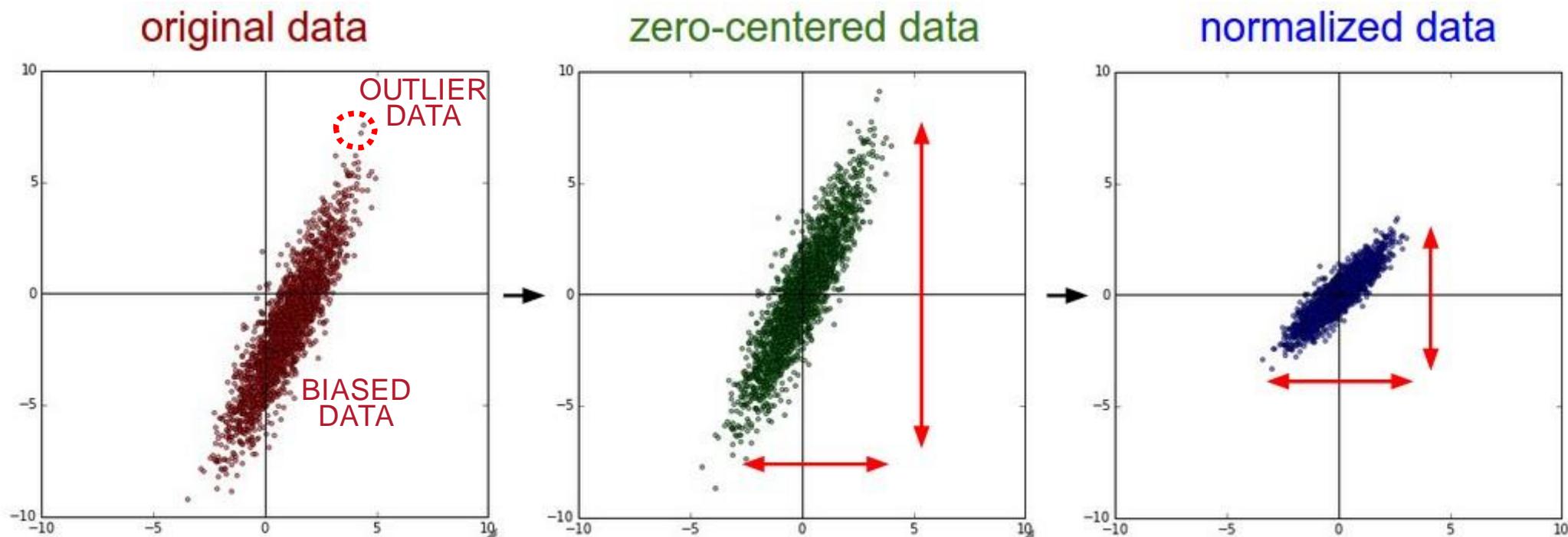
EPOCH & ITERATION & BATCH



BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

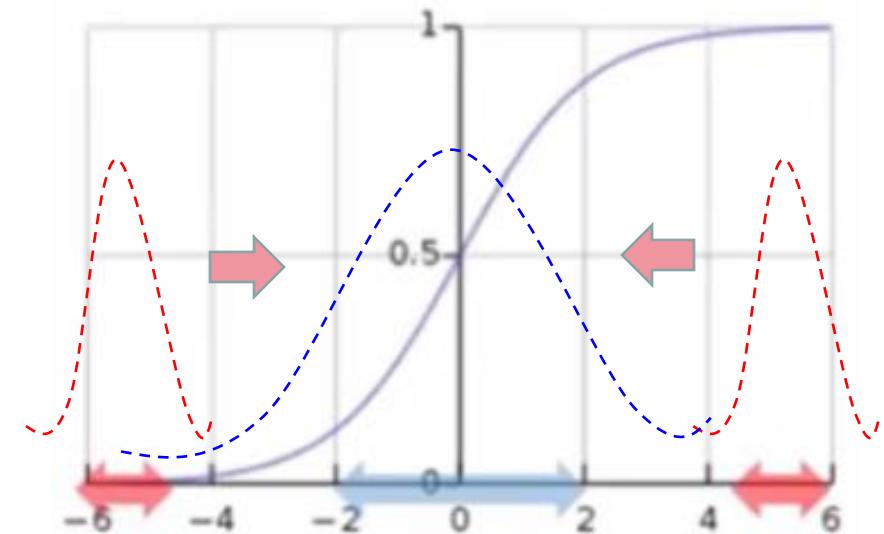
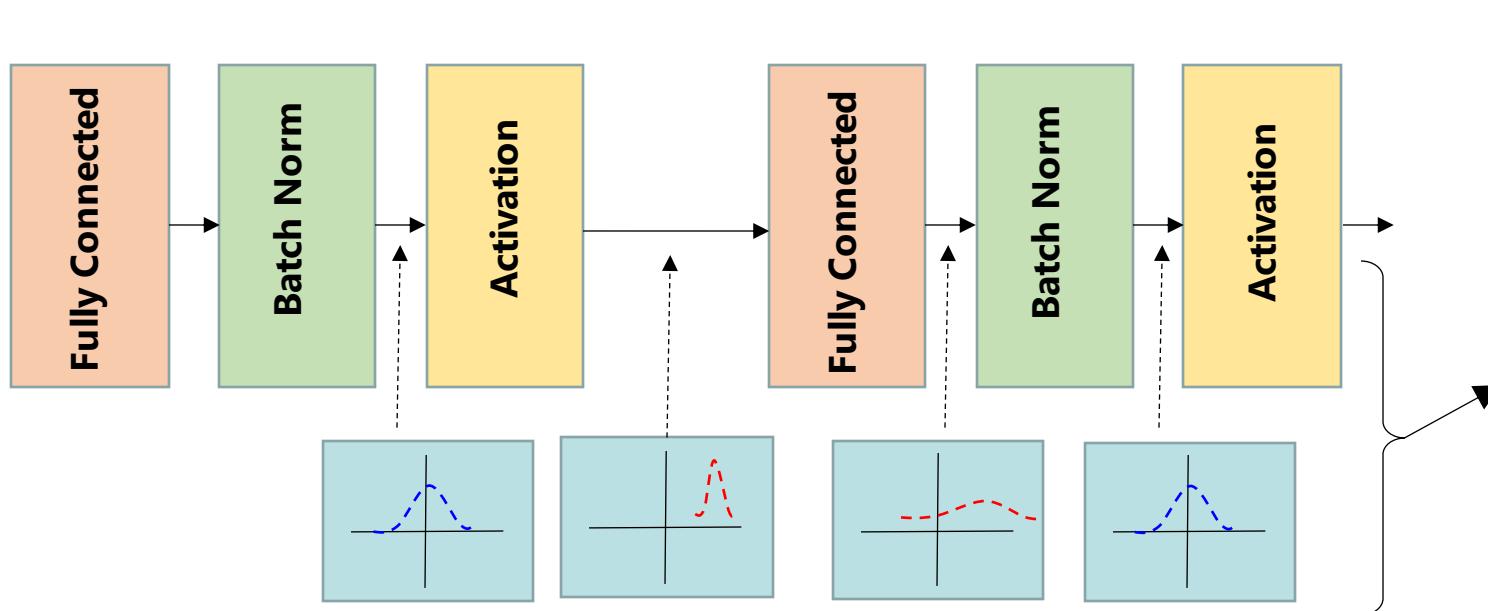
Method to help the network learn better



BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better

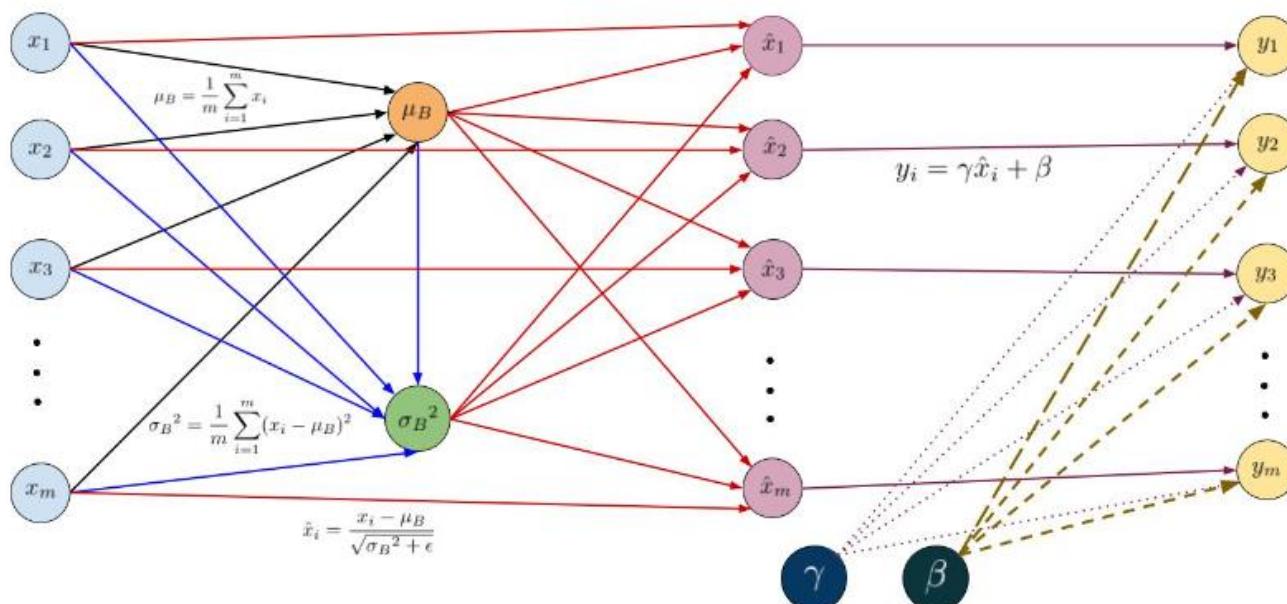


γ, β

BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\begin{aligned}\mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i && // \text{mini-batch mean} \\ \sigma_{\mathcal{B}}^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 && // \text{mini-batch variance} \\ \hat{x}_i &\leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} && // \text{normalize} \\ y_i &\leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) && // \text{scale and shift}\end{aligned}$$

```
1 def batchnorm_forward(x, gamma, beta, eps=1e-5):  
2     N, D = x.shape  
3  
4     sample_mean = x.mean(axis=0)  
5     sample_var = x.var(axis=0)  
6  
7     std = np.sqrt(sample_var + eps)  
8     x_centered = x - sample_mean  
9     x_norm = x_centered / std  
10    out = gamma * x_norm + beta  
11  
12    cache = (x_norm, x_centered, std, gamma)  
13  
14    return out, cache
```

BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better

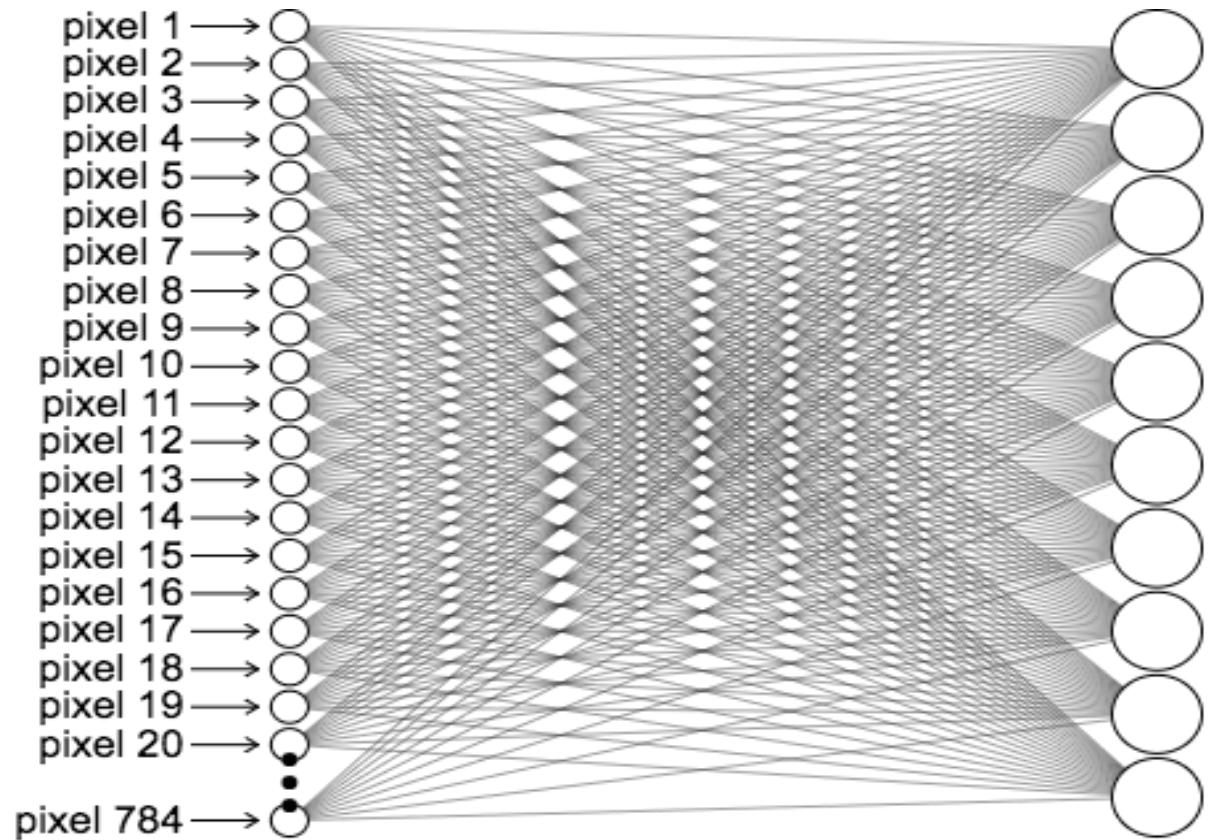
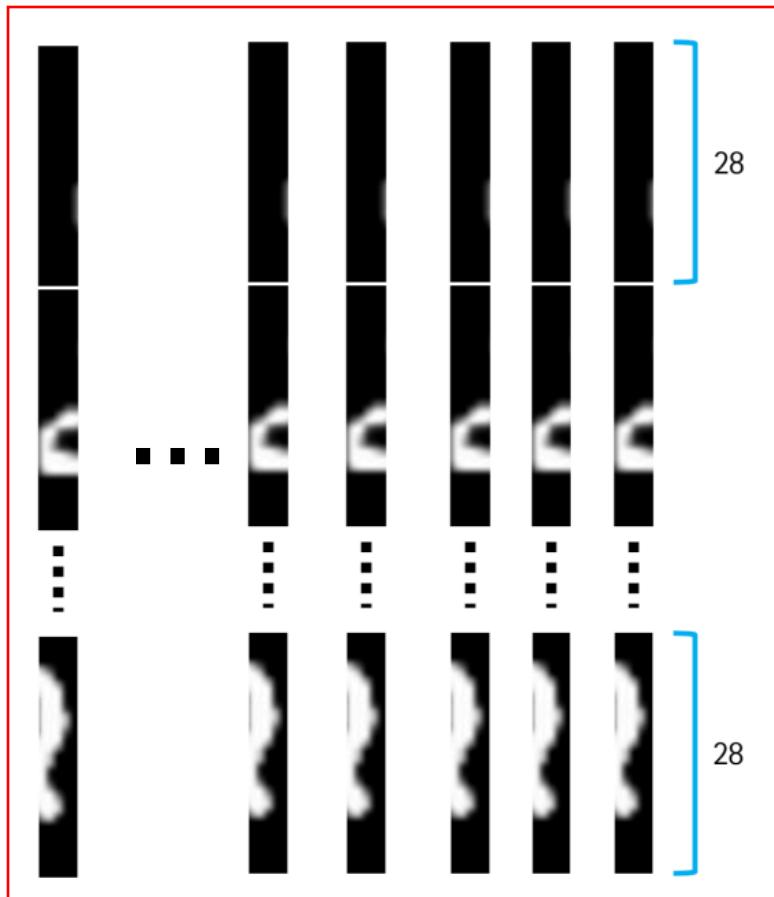
$$\gamma \left(\frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) + \beta \quad \rightarrow$$

scale and shift
→ 분포재조정됨
- 평균: β
- 분산: γ^2

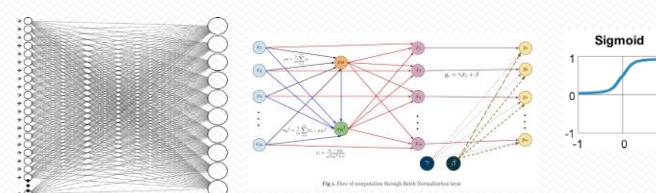
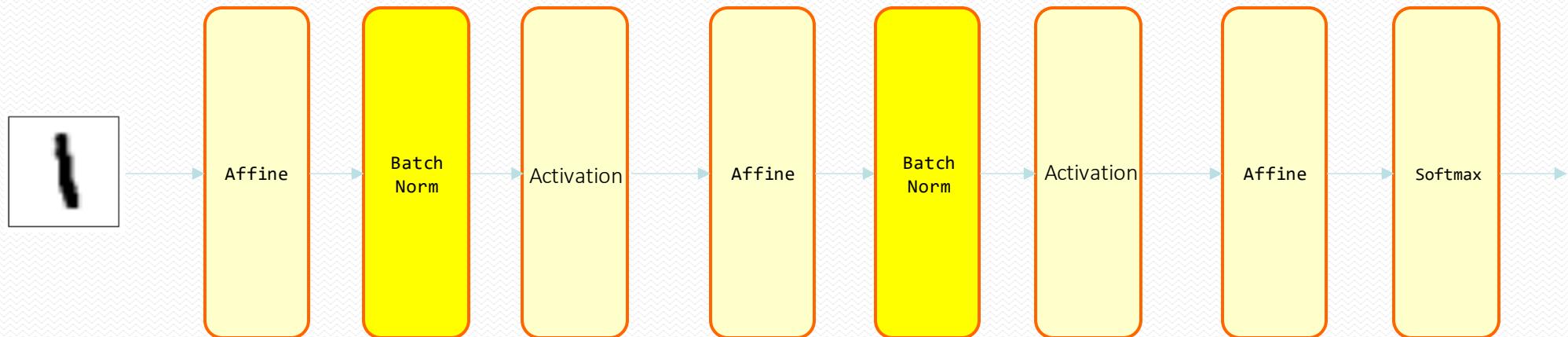
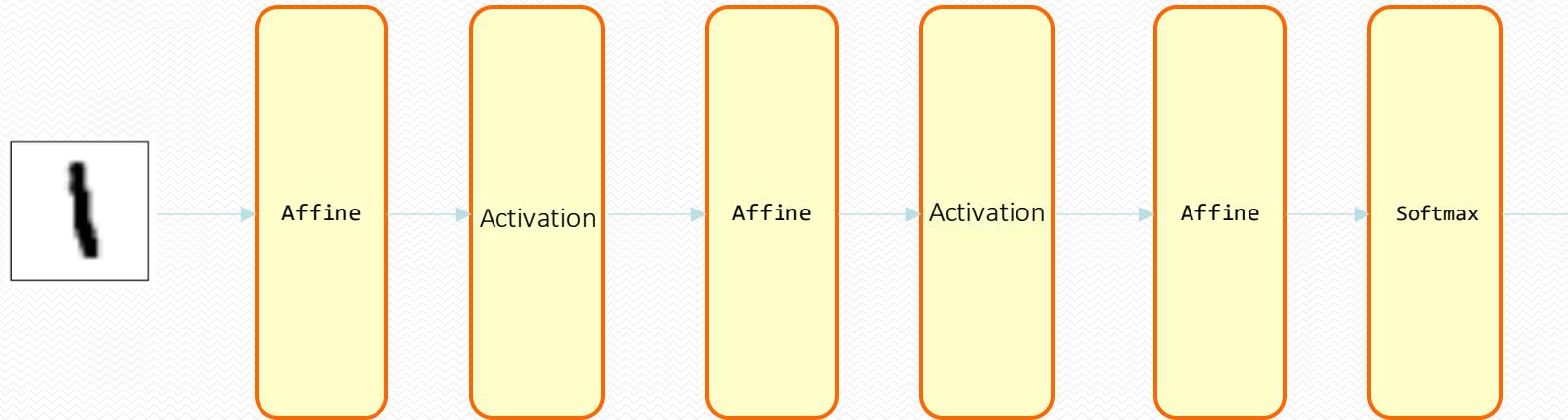
BATCH NORMALIZATION

Batch Normalization

3 4 2 1 9 5 6 2 1 8
8 9 1 2 3 0 0 6 6 4
6 7 0 1 6 3 6 3 7 0
3 7 7 9 4 6 6 1 8 2
2 9 3 4 3 9 8 7 2 5
1 5 9 8 3 6 5 7 2 3
9 3 1 9 1 5 8 0 8 4
5 6 2 6 8 5 8 8 9 9
3 7 7 0 9 4 8 5 4 3
7 9 6 4 1 0 6 9 2 3



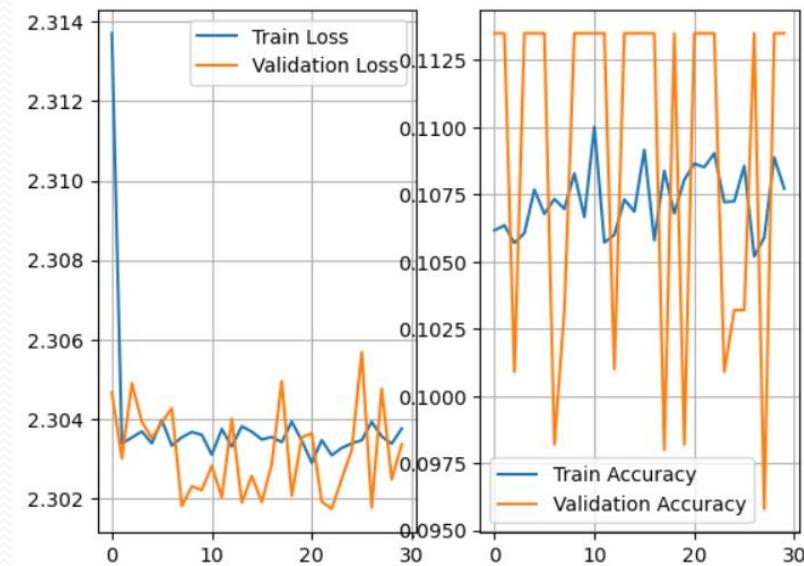
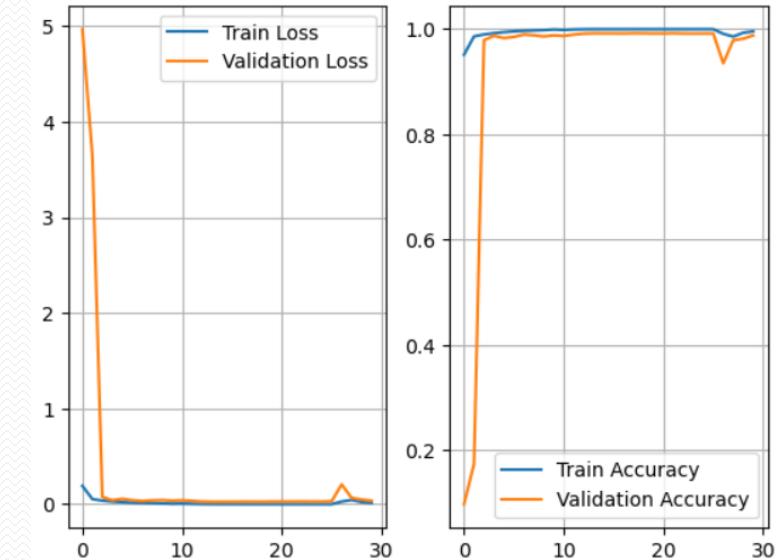
BATCH NORMALIZATION



배치정규화 알고리즘

```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```

```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 # model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 # model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 # model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 # model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 # model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 # model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```



배치정규화 알고리즘

Layer (type)	Output Shape	Param #
<hr/>		
conv2d (Conv2D)	(None, 27, 27, 32)	160
batch_normalization (Batch Normalization)	(None, 27, 27, 32)	128
conv2d_1 (Conv2D)	(None, 26, 26, 64)	8256
batch_normalization_1 (BatchNormalization)	(None, 26, 26, 64)	256
conv2d_2 (Conv2D)	(None, 13, 13, 128)	32896
batch_normalization_2 (BatchNormalization)	(None, 13, 13, 128)	512
conv2d_3 (Conv2D)	(None, 12, 12, 32)	16416
batch_normalization_3 (BatchNormalization)	(None, 12, 12, 32)	128
conv2d_4 (Conv2D)	(None, 11, 11, 64)	8256
batch_normalization_4 (BatchNormalization)	(None, 11, 11, 64)	256
conv2d_5 (Conv2D)	(None, 5, 5, 128)	32896
batch_normalization_5 (BatchNormalization)	(None, 5, 5, 128)	512
flatten (Flatten)	(None, 3200)	0
dense (Dense)	(None, 128)	409728
dense_1 (Dense)	(None, 10)	1290
<hr/>		

Total params: 511690 (1.95 MB)
Trainable params: 510794 (1.95 MB)
Non-trainable params: 896 (3.50 KB)

None

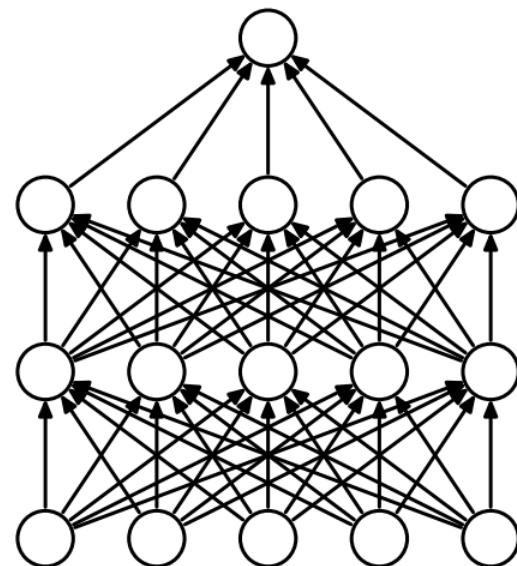
```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```

DROPOUT

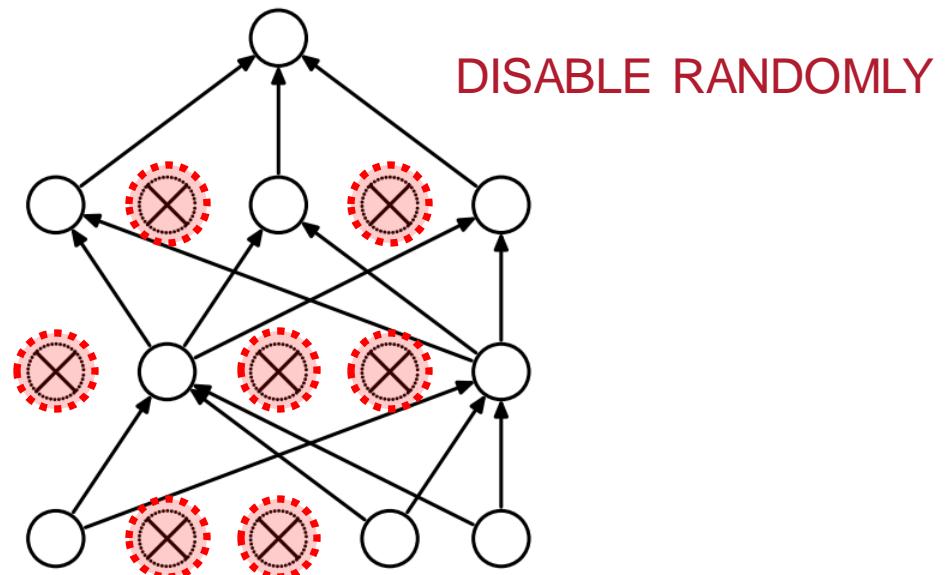
RANDOMLY SET SOME NEURONS TO ZERO IN THE FORWARD PASS

Voting effect → makes the result not biased

Anti co-adaptation → creates a robust network



(a) Standard Neural Net

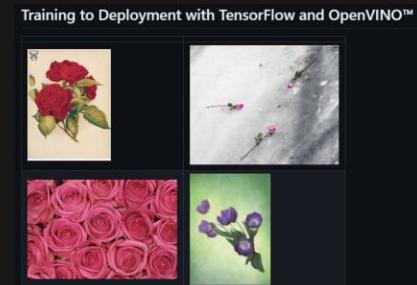


(b) After applying dropout.

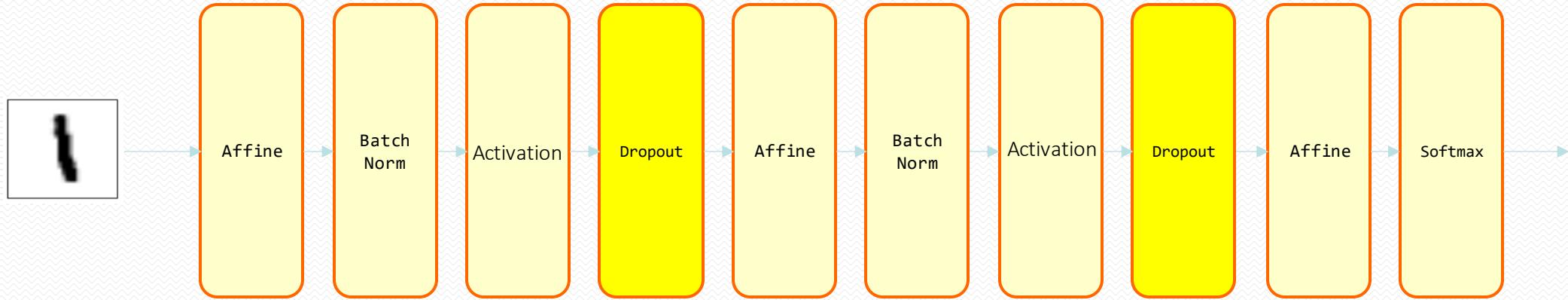
SIMPLE CNN EXAMPLE WITH KERAS

Layer (type)	Output Shape	Param #
<hr/>		
rescaling_1 (Rescaling)	(None, 180, 180, 3)	0
conv2d (Conv2D)	(None, 180, 180, 16)	448
max_pooling2d (MaxPooling2D)	(None, 90, 90, 16)	0
conv2d_1 (Conv2D)	(None, 90, 90, 32)	4640
max_pooling2d_1 (MaxPooling2D)	(None, 45, 45, 32)	0
conv2d_2 (Conv2D)	(None, 45, 45, 64)	18496
max_pooling2d_2 (MaxPooling2D)	(None, 22, 22, 64)	0
flatten (Flatten)	(None, 30976)	0
dense (Dense)	(None, 128)	3965056
dense_1 (Dense)	(None, 5)	645
<hr/>		
Total params: 3,989,285		
Trainable params: 3,989,285		
Non-trainable params: 0		

```
# 모델 만들기
model = tf.keras.Sequential()
model.add(
    tf.keras.layers.Rescaling(
        1./255, input_shape=(IMAGE_WIDTH, IMAGE_HEIGHT, 3)))
model.add(
    tf.keras.layers.Conv2D(16, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Conv2D(32, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Conv2D(64, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Flatten())
model.add(
    tf.keras.layers.Dropout(0.2))
model.add(
    tf.keras.layers.Dense(128, activation="relu"))
model.add(
    tf.keras.layers.Dense(5))
model.summary()
```



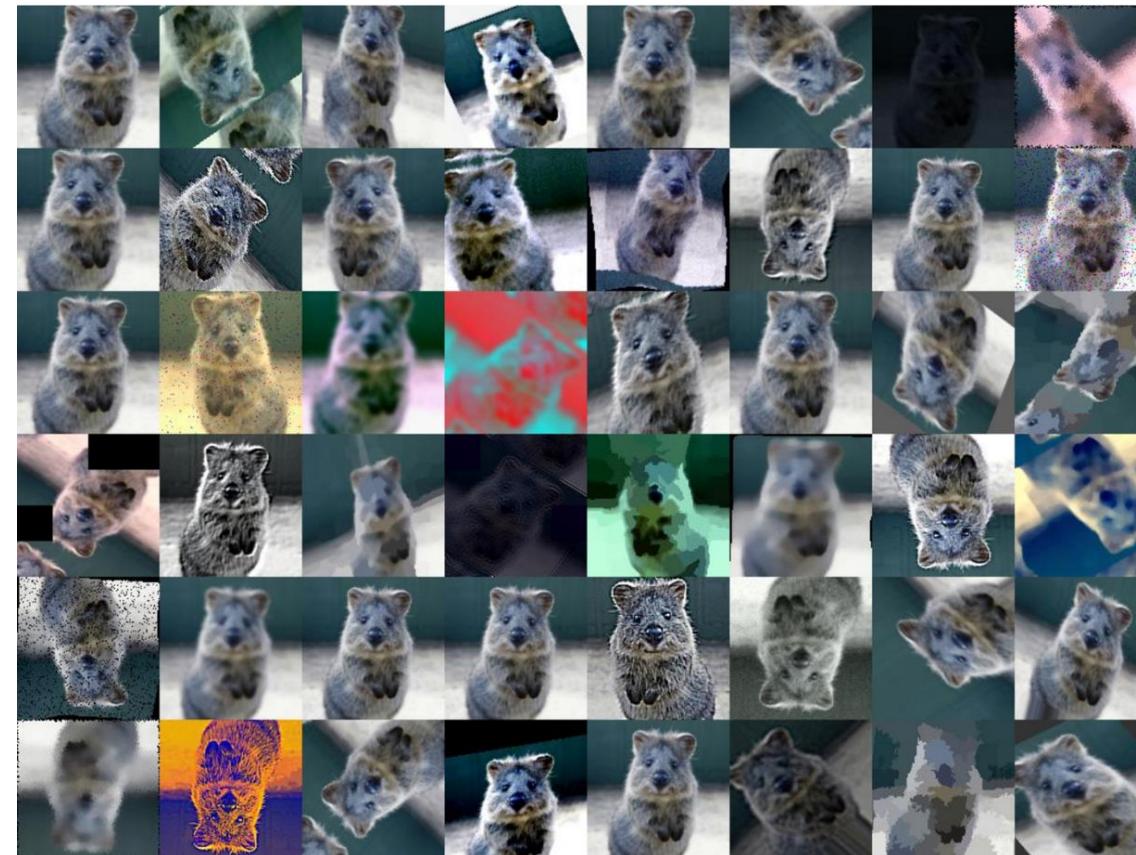
드롭아웃

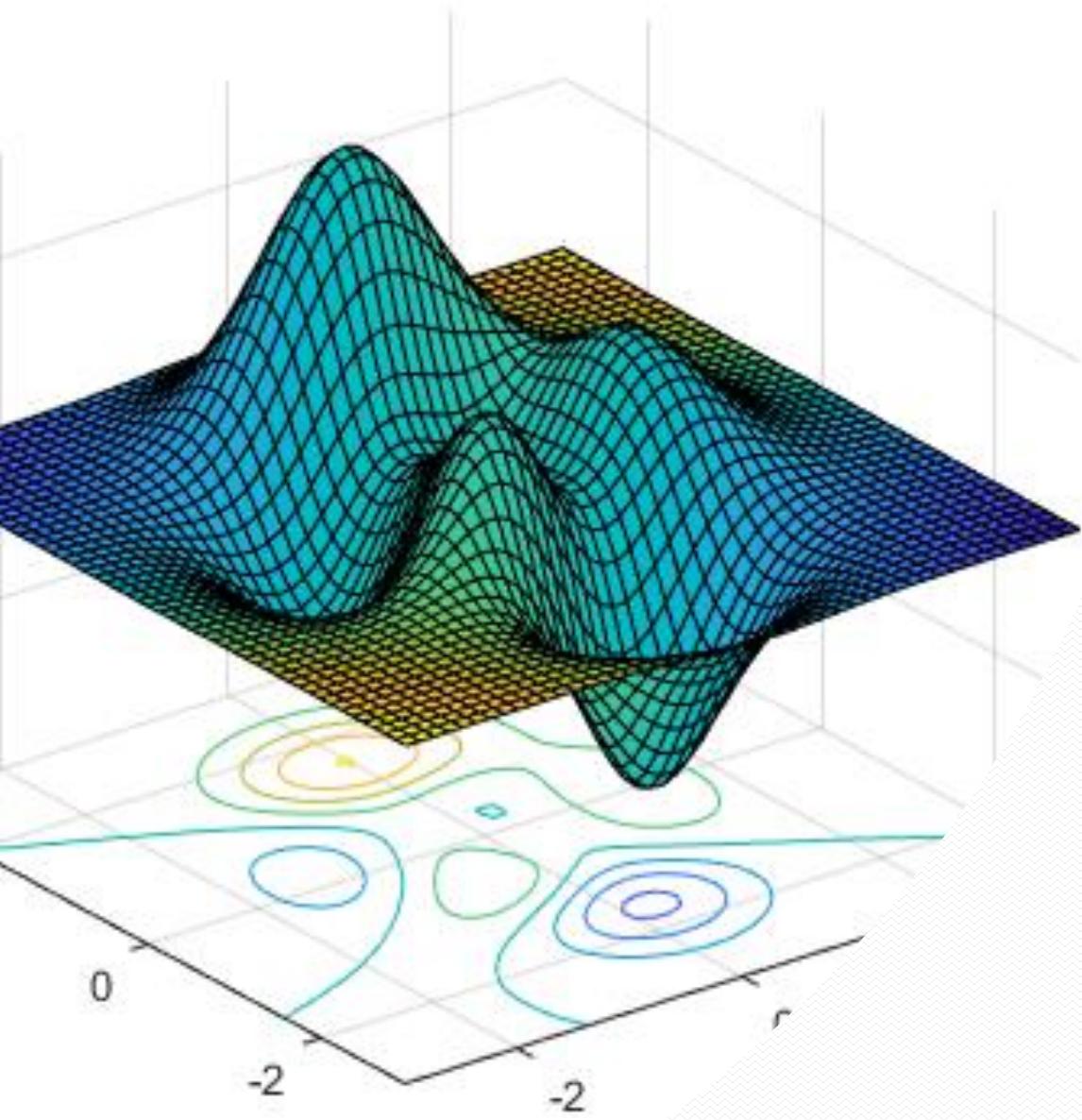


GENERALIZATION / REGULARIZATION

DATA AUGMENTATION

Strategy for improve CNN performance that enables user to significantly increase the diversity of data available for training models, without collecting new data





Deep-dive

Optimizer

Optimizer: Gradient descent algorithm

```
1 # Let's build the CNN model
2
3 cnn = Sequential()
4
5 # Convolution
6 cnn.add(Conv2D(32, (3, 3), activation="relu", input_shape=(64, 64, 3)))
7
8 # Pooling
9 cnn.add(MaxPooling2D(pool_size = (2, 2)))
10
11 # 2nd Convolution
12 cnn.add(Conv2D(32, (3, 3), activation="relu"))
13
14 # 2nd Pooling layer
15 cnn.add(MaxPooling2D(pool_size = (2, 2)))
16
17 # Flatten the layer
18 cnn.add(Flatten())
19
20 # Fully Connected Layers
21 cnn.add(Dense(activation = 'relu', units = 128))
22 cnn.add(Dense(activation = 'sigmoid', units = 1))
23
24 # Compile the Neural network
25 cnn.compile(optimizer = 'adam', loss = 'binary_crossentropy', metrics = ['accuracy'])
26
27
28
29
30
```

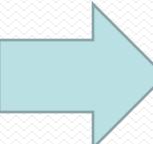
Gradient descent algorithm

Solve the following minimization problem : for $x \in \mathbf{R}^n$,

$$\min_x f(x)$$

To solve the problem, we use descent method framework : For given initial $x^{(0)}$,

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \Delta x^{(k)}$$

- | | |
|---|--|
| 1. $\alpha^{(k)}$: Learning Rate / Step Size
2. $\Delta x^{(k)}$: Search Direction | 
$\alpha(\text{constant})$
$-\nabla f(x^k)$ |
|---|--|

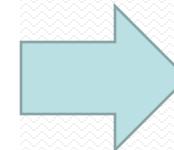
- 비용 함수 최소화
- Gradient descent는 주로 최적화 문제에 사용된다.
- 주어진 비용 함수, 비용 (W, b)에 대해 비용을 최소화하기 위해 W, b 를 찾는다.
- 일반적인 함수 : 비용 (w_1, w_2, \dots)에 적용 가능

Gradient descent algorithm

Adaptive Learning rate scheme

- Learning rate가 작으면 안정적이지만 너무 느리고
- Learning rate가 크면 빠르지만 불안정
- 두가지 상황에서 단점을 빼고 장점만 가져가는 방법은?

1. $\alpha^{(k)}$: Learning Rate / Step Size
2. $\Delta x^{(k)}$: Search Direction



$$\begin{aligned} &??? \alpha \\ &??? -\nabla f(x^k) \end{aligned}$$

Gradient descent optimization Algo.

BGD (Batch Gradient Descent)

- FULL TRAIN DATASET, REQUIRES A LOT OF COMPUTATION

SGD (Stochastic Gradient Descent)

- USE Some of training sample, calculate gradient

Momentum

- Method that helps accelerate SGD in the relevant direction and dampens oscillations by adding update vector of the past time step to the current update vector



Image 2: SGD without momentum



Image 3: SGD with momentum

Gradient descent optimization Algo.

Adagrad (Adaptive Gradient)

- Adapts learning rate to the parameters, performing smaller updates(i.e. low learning rates) for parameters associated with frequently occurring features, and larger updates (i.e. high learning rates) for parameters associated with infrequent features.

RMSProp

- Need to resolve Adagrad's radically diminishing learning rates. RMSprop divides the learning rate by an exponentially decaying average of squared gradients

Adam (Adaptive Moment Estimation)

- Another approach to compute adaptive learning rates for each parameter. In addition to storing an exponentially decaying average of past squared gradients like Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients similar to momentum

Gradient descent algorithm

Momentum : 1945

- 방식은 진동하는 현상을 겪을 때 이를 해결하도록 해 줌
- 이전 gradient들도 계산에 포함해서 현재 파라미터를 업데이트
- 이전 gradient들을 모두 동일한 비율로 포함시키지는 않고 비율을 감소시켜줌

$$v_{k+1} = \alpha v_k - \epsilon \nabla f(x_k)$$

$$x_{k+1} = x_k + v_{k+1}$$

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while stopping criterion not met do

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

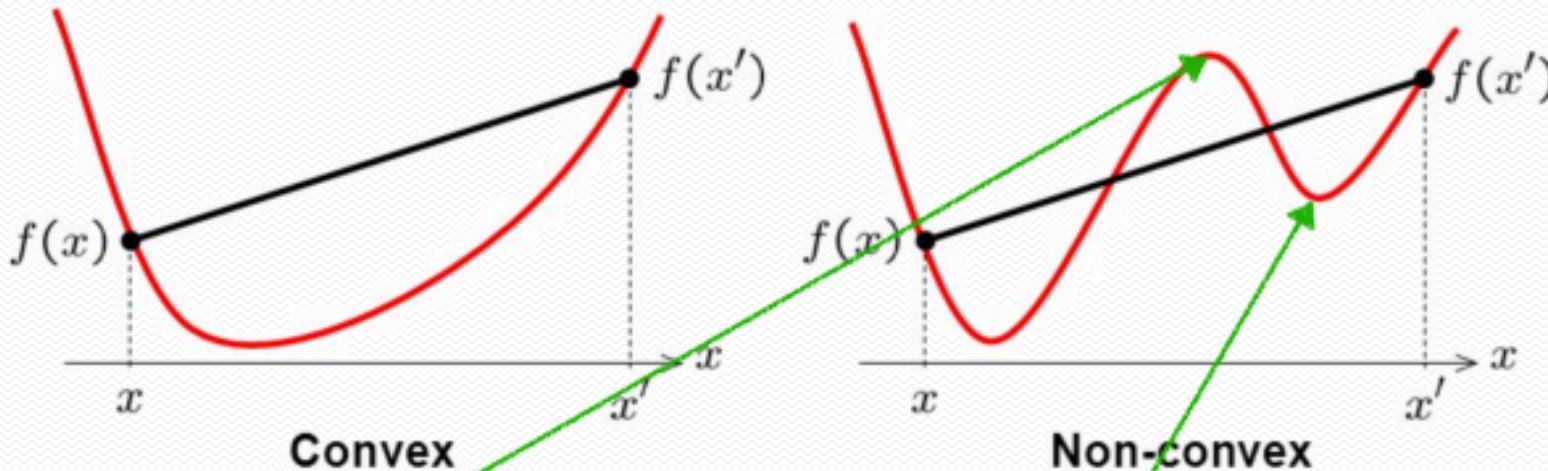
 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$.

 Apply update: $\theta \leftarrow \theta + v$.

end while

Gradient descent algorithm

Momentum : 1945



- 미분값이 0이라고 무조건 최솟값은 아닙니다.
- **최댓값**일 수도 있습니다.
- 최솟값은 아니지만 미분 값이 0인 곳을 **Local Minimum**이라고 합니다.

Gradient descent algorithm

Adagrad : 2011

- 각 weight마다 다른 learning rate를 줌
- 변화한 누적거리를 계속 저장하여 learning rate에 반영
- 변화한 누적 거리가 큰 weight는 작은 learning rate를 반영
- 변화한 누적 거리가 작은 weight는 큰 learning rate를 반영

$$r_{k+1} = r_k + \nabla f(x_k) \odot \nabla f(x_k)$$
$$x_{k+1} = x_k - \frac{\epsilon}{\delta + \sqrt{r_{k+1}}} \odot \nabla f(x_k)$$

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $r = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $r \leftarrow r + \mathbf{g} \odot \mathbf{g}$.

 Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$.

end while

Gradient descent algorithm

RMSprop : 2012

- Adagrad와 비슷한 알고리즘
- 아주 오래전 궤적의 영향을 줄이고 가까운 시간의 궤적의 영향을 높이기 위해 decay를 추가

$$r_{k+1} = \rho r_k + (1 - \rho) \nabla f(x_k) \odot \nabla f(x_k)$$
$$x_{k+1} = x_k - \frac{\epsilon}{\sqrt{\delta + r_{k+1}}} \odot \nabla f(x_k)$$

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers

Initialize accumulation variables $r = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $r \leftarrow \rho r + (1 - \rho) \mathbf{g} \odot \mathbf{g}$.

 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + r}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + r}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$.

end while

Gradient descent algorithm

Adam : 2014

- Momentum과 RMSprop를 Hybrid한 방법
- Hyper-parameter 설정에 가장 robust하다고 알려짐 (Bengio, 2016]

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization (Suggested default:
 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $s = \mathbf{0}$, $r = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with
 corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

 Update biased first moment estimate: $\hat{s} \leftarrow \rho_1 s + (1 - \rho_1) \mathbf{g}$

 Update biased second moment estimate: $\hat{r} \leftarrow \rho_2 r + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

 Correct bias in first moment: $\hat{s} \leftarrow \frac{\hat{s}}{1 - \rho_1^t}$

 Correct bias in second moment: $\hat{r} \leftarrow \frac{\hat{r}}{1 - \rho_2^t}$

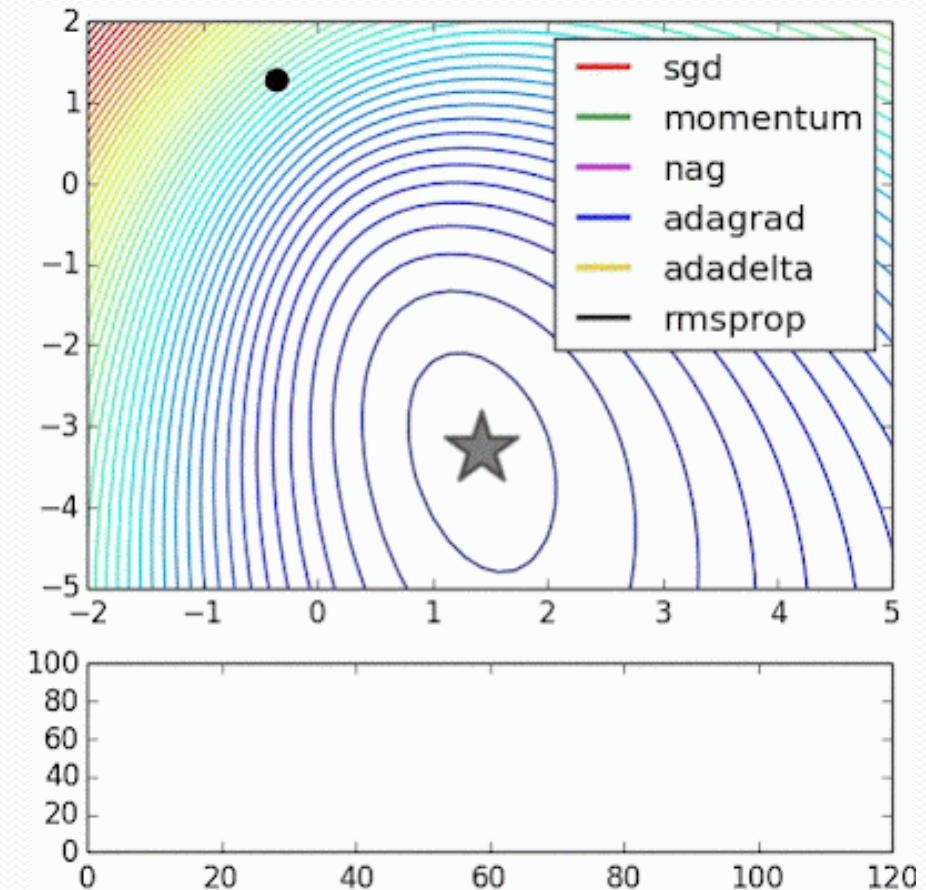
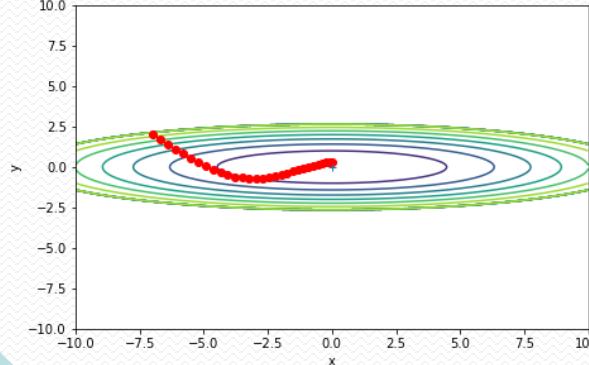
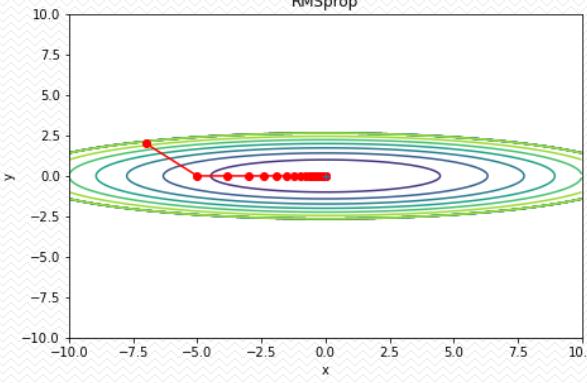
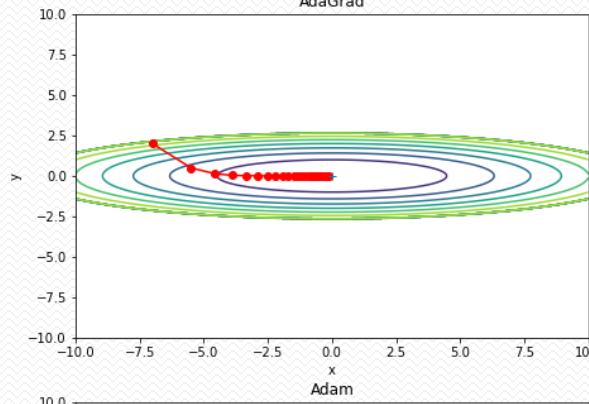
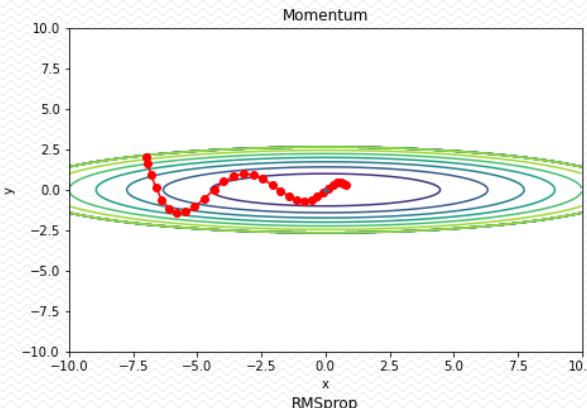
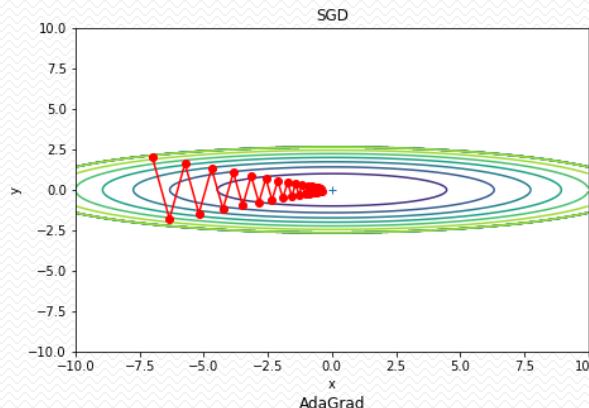
 Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}$ (operations applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

$$\begin{aligned}s_{k+1} &= \rho_1 s_k + (1 - \rho_1) \nabla f(x_k) \\r_{k+1} &= \rho_2 r_k + (1 - \rho_2) \nabla f(x_k) \odot \nabla f(x_k) \\\hat{s}_{k+1} &= \frac{s_{k+1}}{1 - \rho_1^t} \\\hat{r}_{k+1} &= \frac{r_{k+1}}{1 - \rho_2^t} \\x_{k+1} &= x_k - \frac{\epsilon}{\sqrt{\hat{r}_{k+1}} + \delta} \hat{s}_{k+1}\end{aligned}$$

Behavior & Performance



CALCULATE THE OUTPUT DATA / SOFTMAX

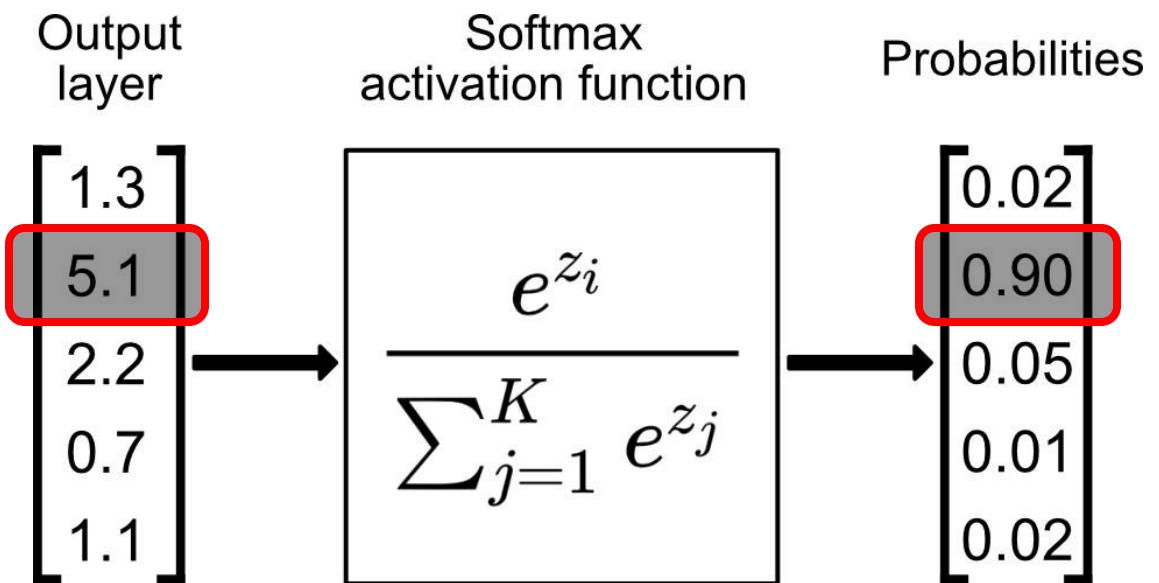
CALCULATE THE CONFIDENCE SCORE (Probabilities)

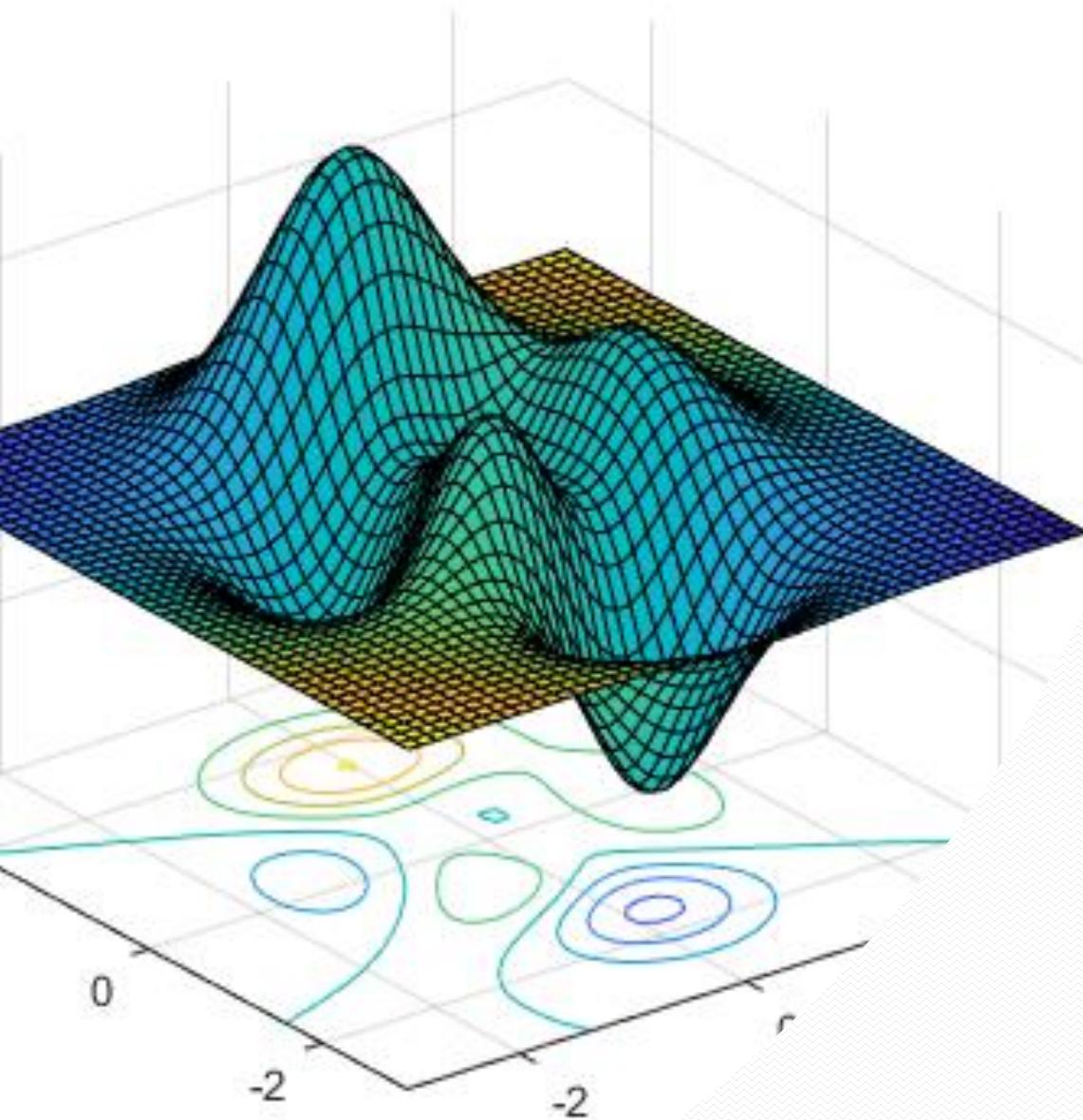
- Suitable for classification problems
multi-class classification
- Sum of scores will be 1
- Highest score is the prediction class

Hypothesis :
 $\hat{y} = wx + b$

Parameters:
 w, b

Cost(Loss) Function
 $C_{\text{cost}}(w, b) = \frac{1}{2}(\hat{y} - y)^2$





Back Propagation

MLP TRAINING & INFERENCE

FEEDFORWARD

Predict

BACKPROPAGATE

Weight Update

FEEDFORWARD

Predict

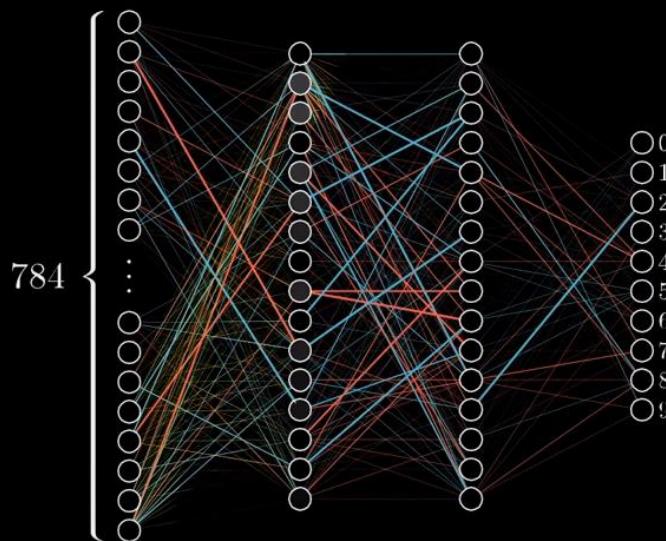
BACKPROPAGATE

Weight Update

...

Training in
progress...

 → 9



chain rule

$$x \rightarrow x^2 \rightarrow x^2 + 1 \rightarrow (x^2 + 1)^2$$
$$f(x) = (x^2 + 1)^2$$

$$\frac{d((x^2+1)^2)}{dx} = \frac{d((x^2+1)^2)}{d(x^2+1)} \frac{d(x^2+1)}{dx^2} \frac{dx^2}{dx}$$

$$f=(x^2+1)^2, g=x^2+1, h=x^2$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dh} \frac{dh}{dx}$$

$$\frac{d((x^2+1)^2)}{dx} = \frac{df}{dg} \frac{dg}{dh} \frac{dh}{dx} = 2(x^2+1) \cdot 1 \cdot 2x = 4x(x^2+1)$$

chain rule

$$\frac{\partial f(g(h(x)))}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial x} = f'(g(h(x))) g'(h(x)) h'(x) = (f \circ g \circ h)'$$

$$f(x) = (x^2 + 1)^2$$

$$f(g) = g^2$$

$$g(h) = h + 1$$

$$h(x) = x^2$$

$$f'(g) = 2g$$

$$g'(h) = 1$$

$$h'(x) = 2x$$

Direct method

$$f(x) = (x^2 + 1)^2 = x^4 + 2x^2 + 1$$

$$f'(x) = 4x^3 + 4x = 4x(x^2 + 1)$$

Chain rule

$$f'(x) = 2(x^2 + 1) \cdot 1 \cdot 2x$$

$$= 4x(x^2 + 1)$$

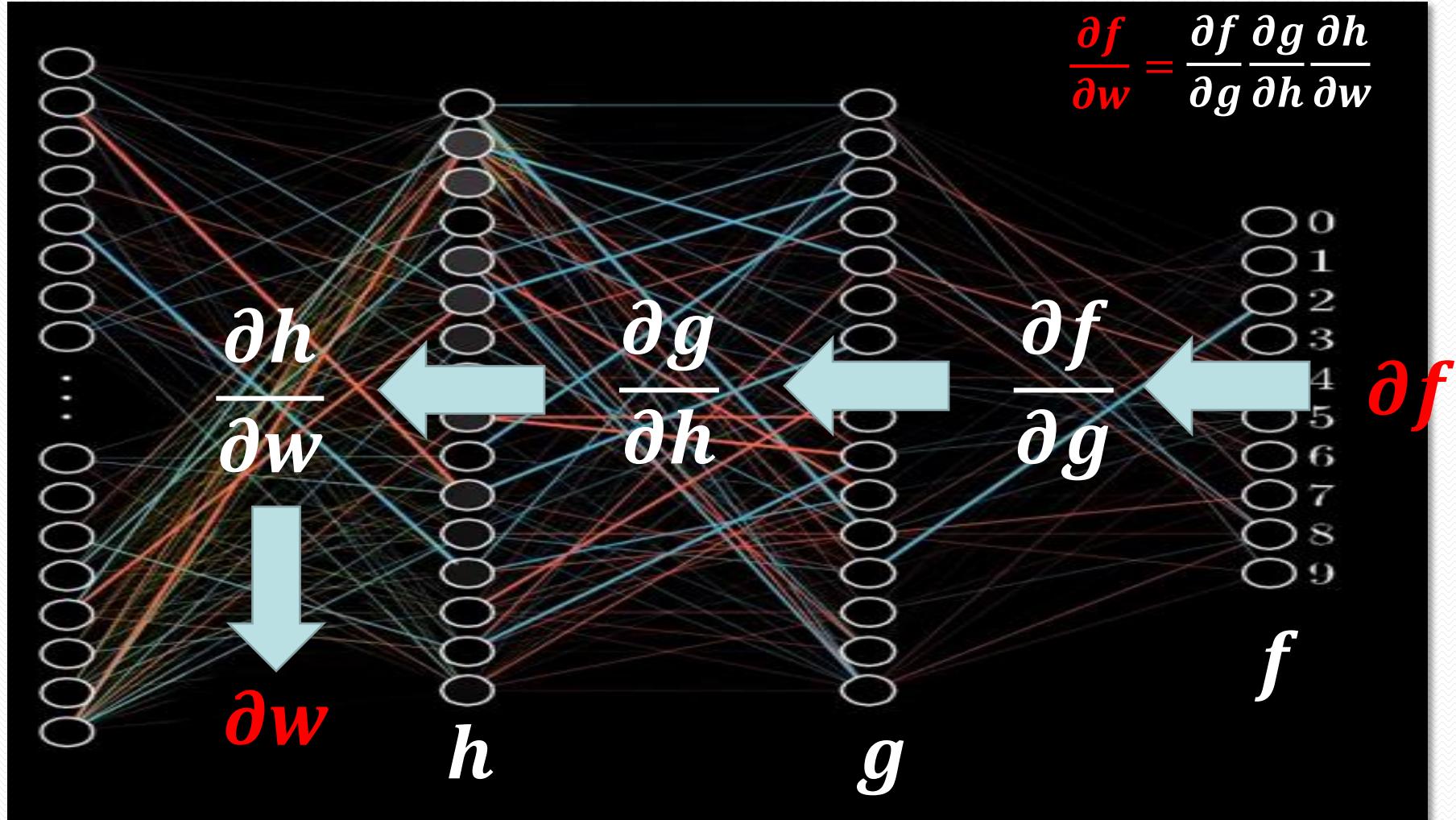
chain rule

$$(f \circ g \circ h)' = \frac{\partial f(g(h))}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial x}$$

- ▶ 체인룰은 변수가 여러 개일 때, 어떤 변수에 대한 다른 변수의 변화율을 알아내기 위해 쓰임
- ▶ Back propagation에서 주로 사용됨.

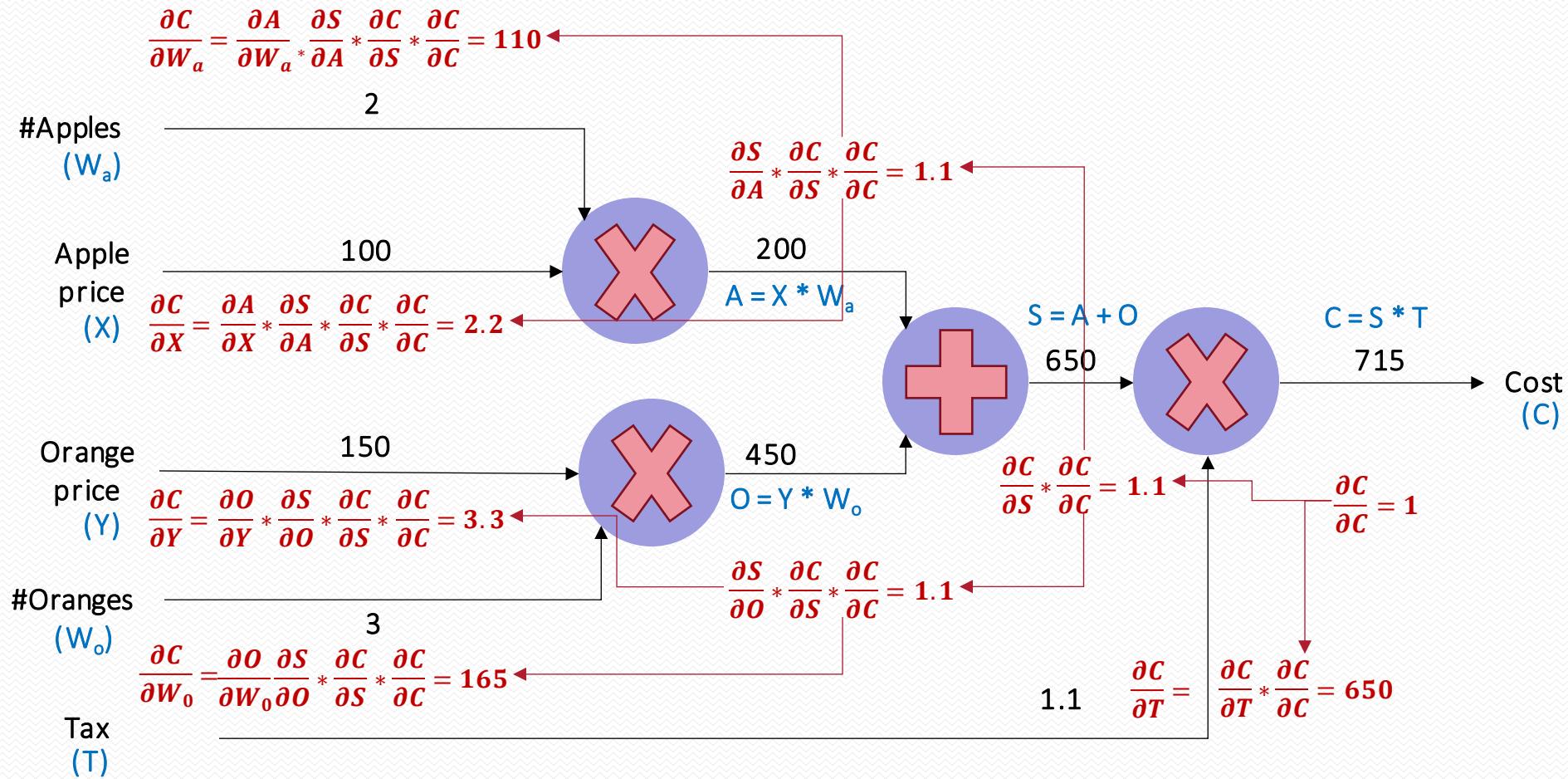
Error Back Propagation

CALCULATE THE ∂w from ∂f



Error Back Propagation

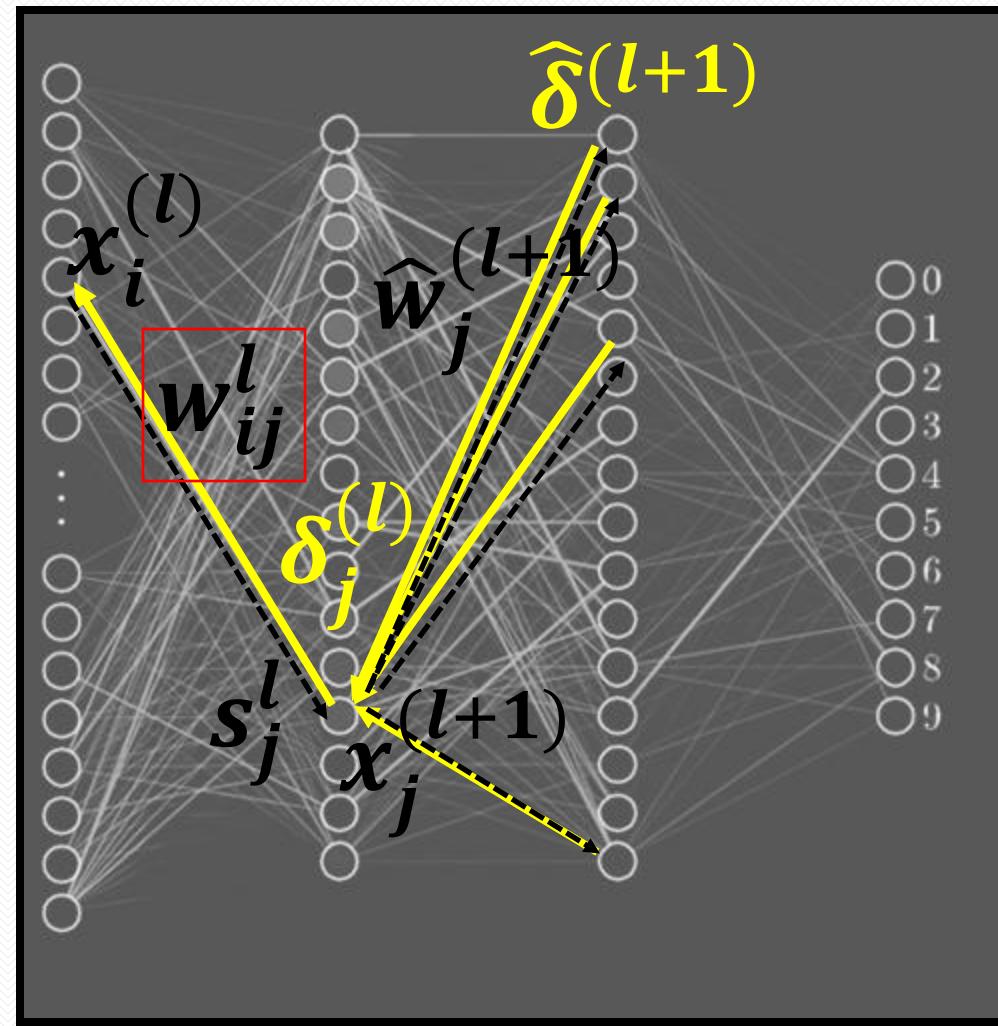
- Contributions of each parameter to final cost can be calculated from partial derivatives.



Back Propagation

$$\frac{\partial f}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} x_i^{(l)}$$

$$\delta_j^{(l)} = (\widehat{\delta}^{(l+1)})^T \widehat{w}_j^{(l+1)} \frac{\partial x_j^{(l+1)}}{\partial s_j^{(l)}}$$





THANK YOU