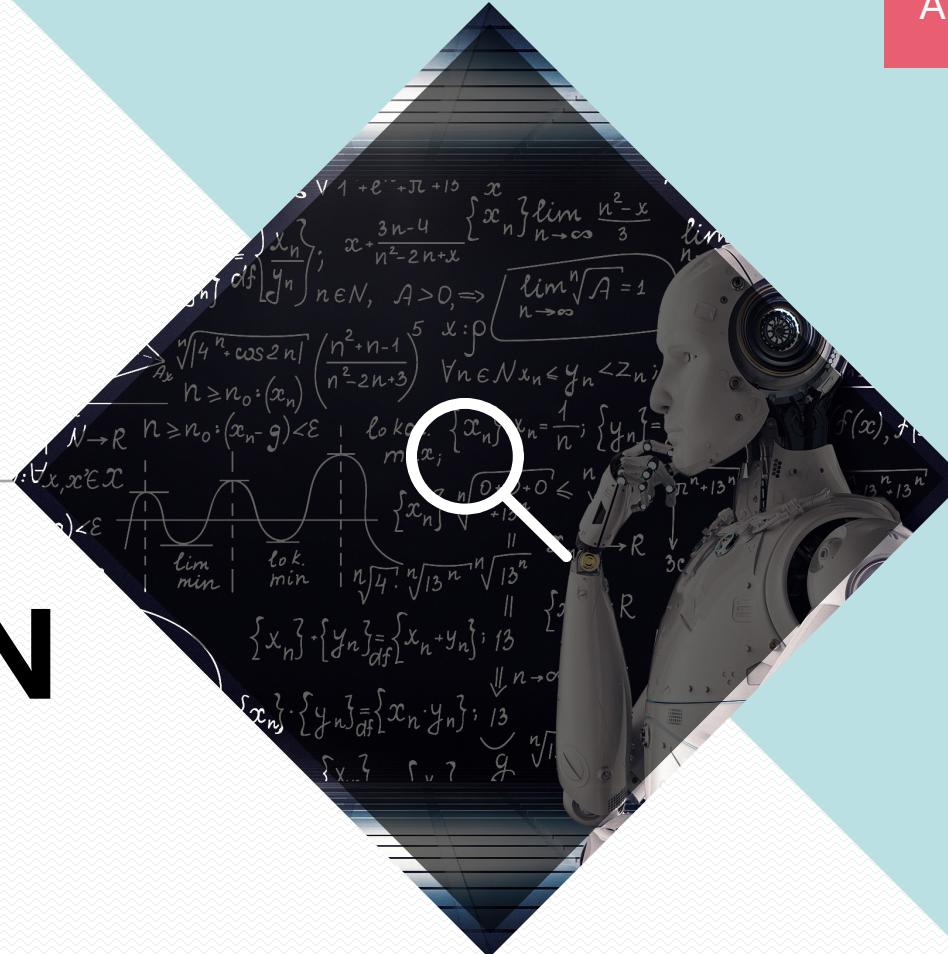


Perceptron & ANN

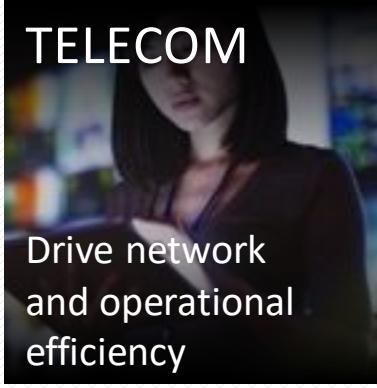
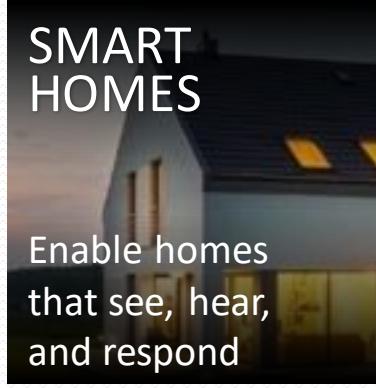
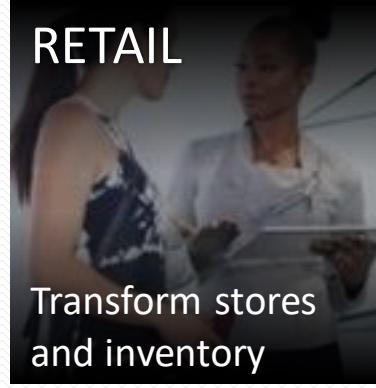
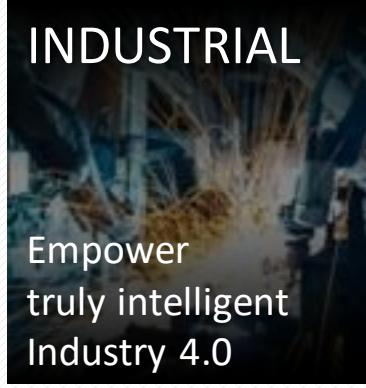
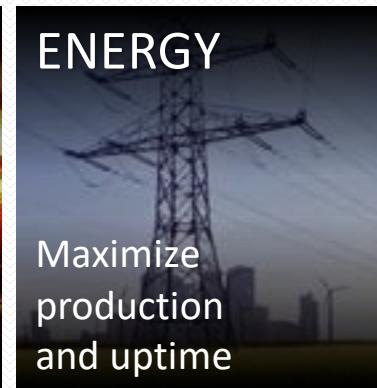
AI/DL introduction & history



AI/DL Momentum



AI is changing every market

[Optimization Notice](#)

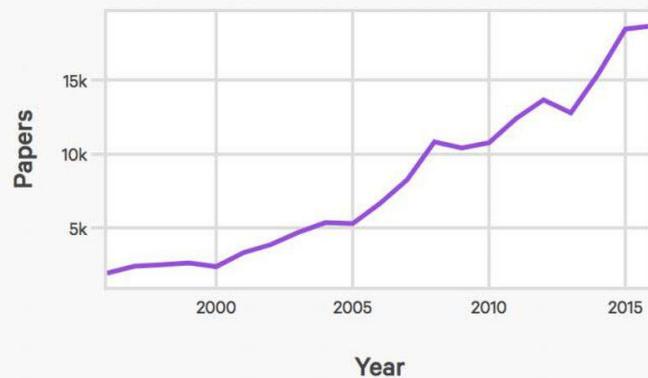
Copyright © 2020, Intel Corporation. All rights reserved.

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

For public use – OK for non-NDA disclosure

AI GROWTH

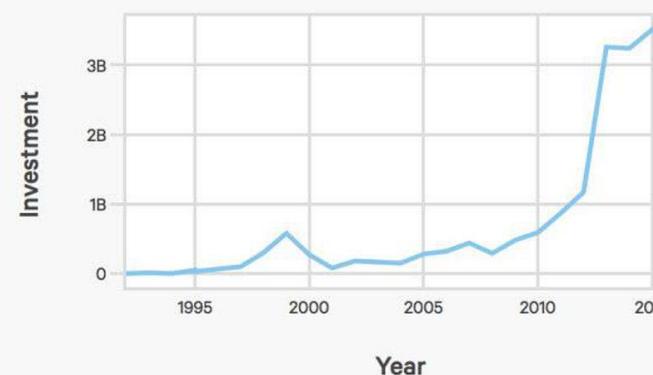
Annually Published AI Papers



Source: Scopus.com

AIINDEX.ORG

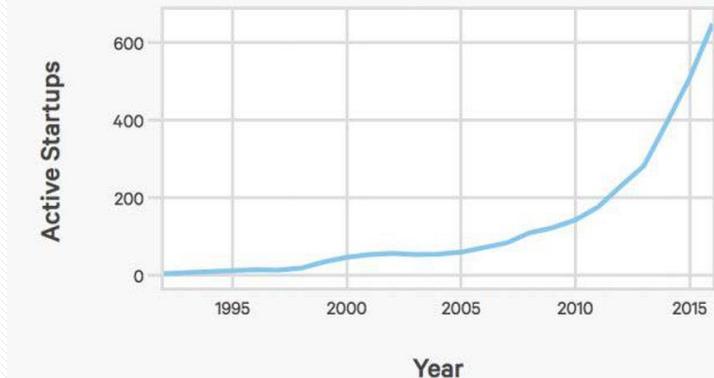
Annual VC Investment in AI Startups



Sources: Crunchbase, VentureSource, Sand Hill Econometrics

AIINDEX.ORG

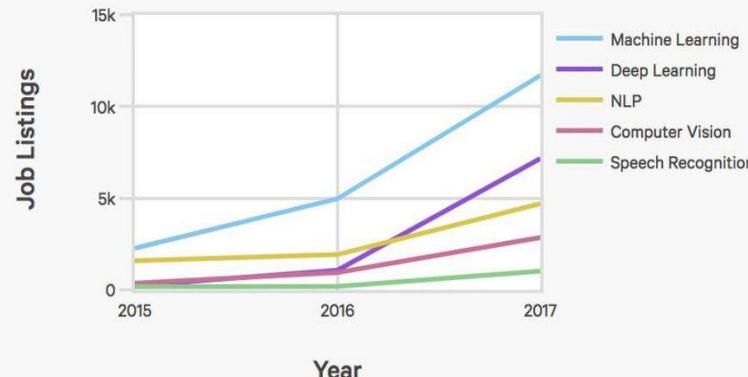
Startups Developing AI Systems



Sources: Crunchbase, VentureSource, Sand Hill Econometrics

AIINDEX.ORG

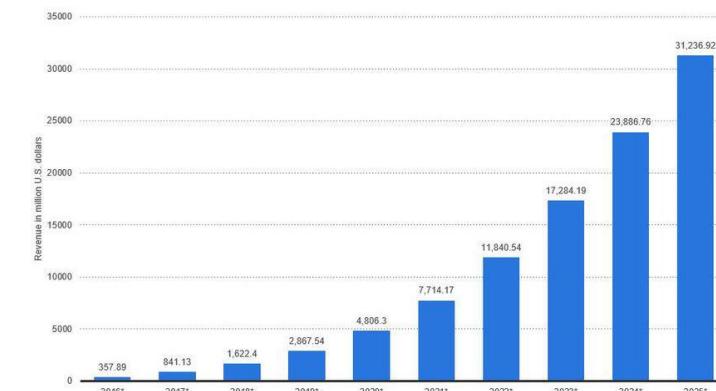
Job Openings, Skills Breakdown (Monster.com)



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)

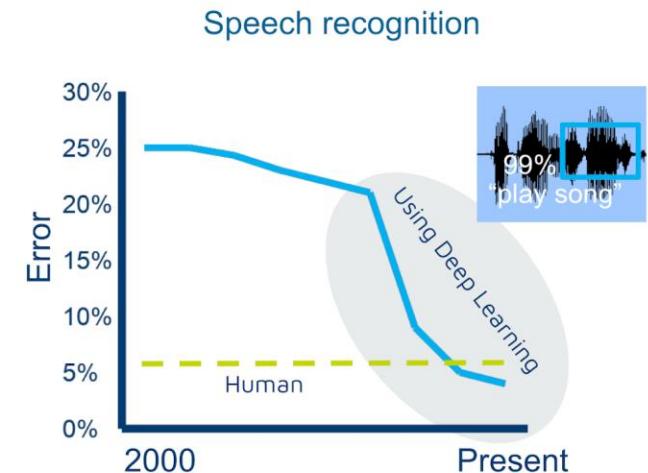
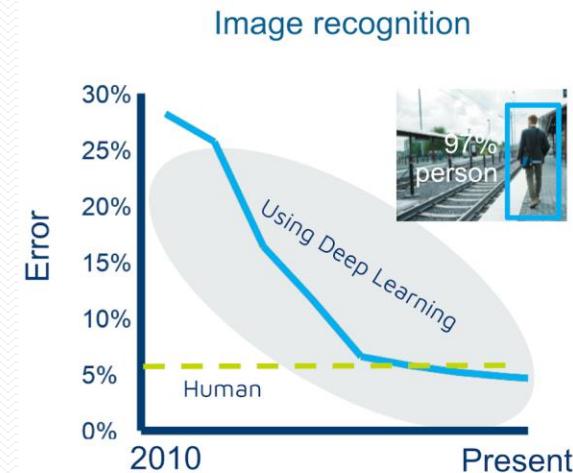


statista

source

WHY NOW?

Bigger Datasets	Better Hardware	Smarter Algorithms
		



AI TRENDS IN MARKET

A.I. Here, There, Everywhere!



RETAIL

Personalized Shopping Experience
Inventory Management
In-store Surveillance
Supply Chain Optimization



INDUSTRIAL

Predictive Maintenance
Machine Vision
Robotics
Improved Quality of Service
Operational Efficiency



TRANSPORT

Enhanced In-cabin Experience
Aircraft Maintenance and Safety
Real-time Tracking of Packages/ Fleet
Autonomous Driving



SMART CITIES

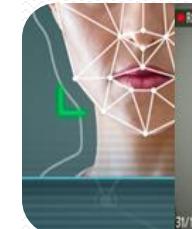
Safety & Security
Resident Engagement
Traffic Management
Object Detection



HEALTHCARE

Enhanced Diagnostics
Drug Discovery
Customized Patient Care
Sensory Aids

AI Products



RETAIL



INDUSTRIAL



TRANSPORT



SMART CITIES



HEALTHCARE



 **LG Energy Solution**  **HYUNDAI**

emart

 Hanwha

 현대중공업 







NAVER

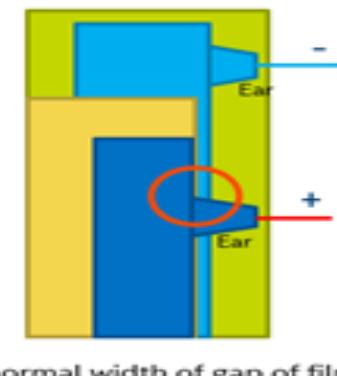
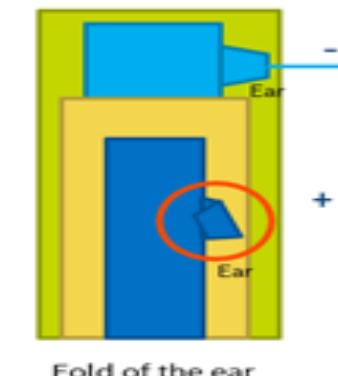
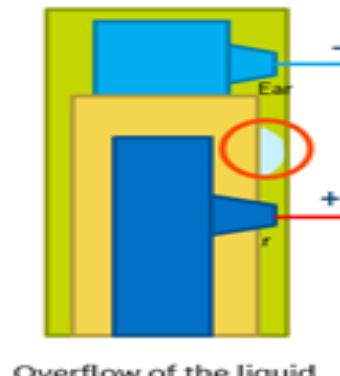
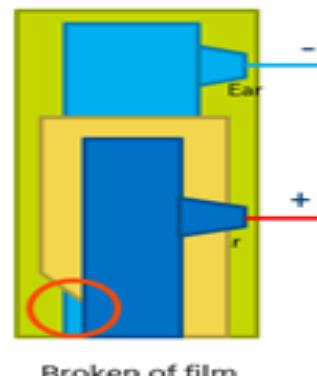
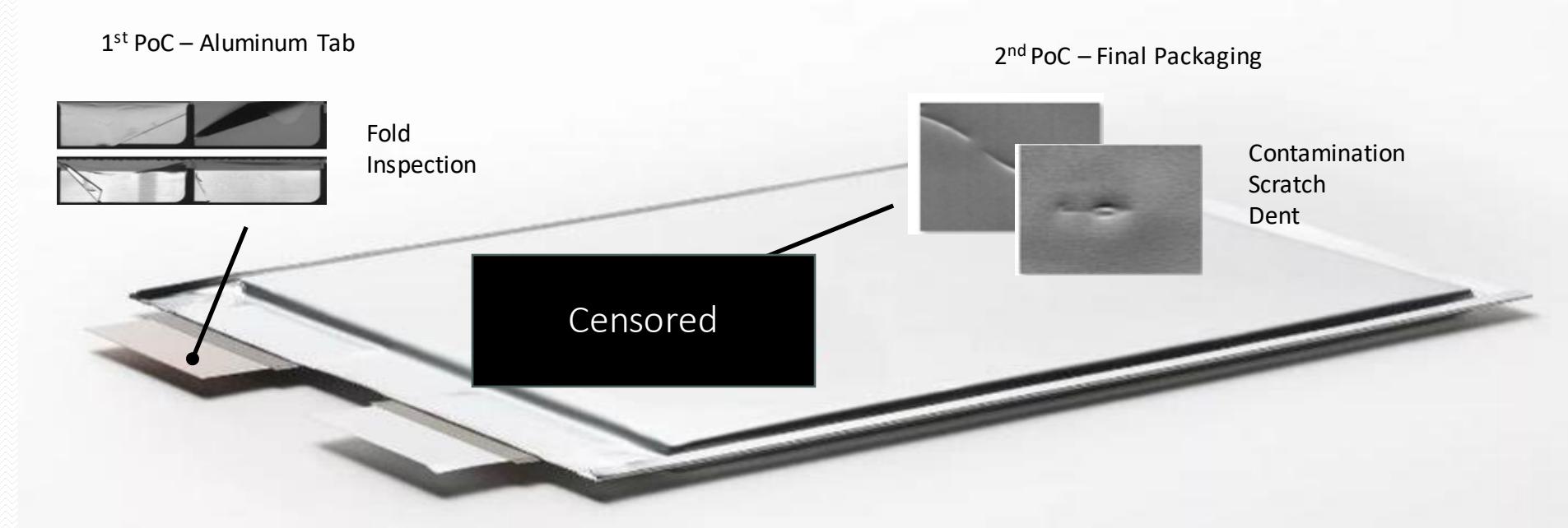
 신한은행



SAMSUNG MEDISON


View the Invisible, Know the Unknown

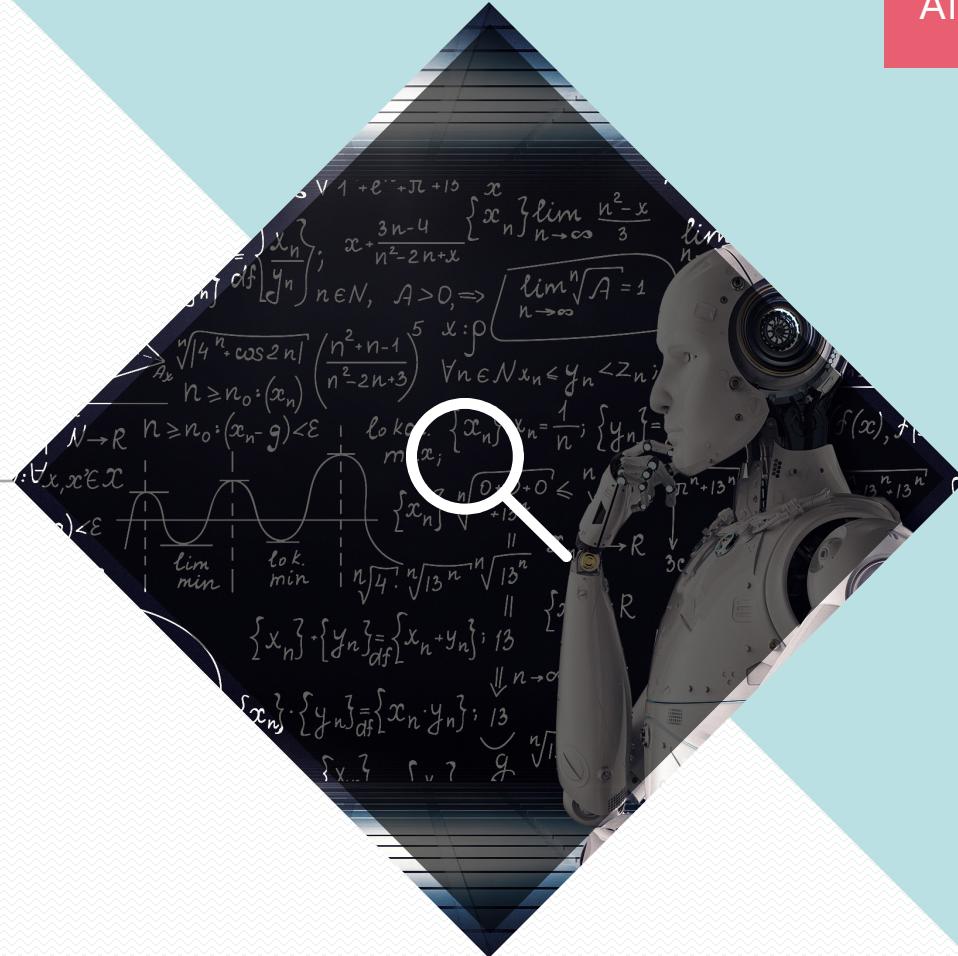
불량검출



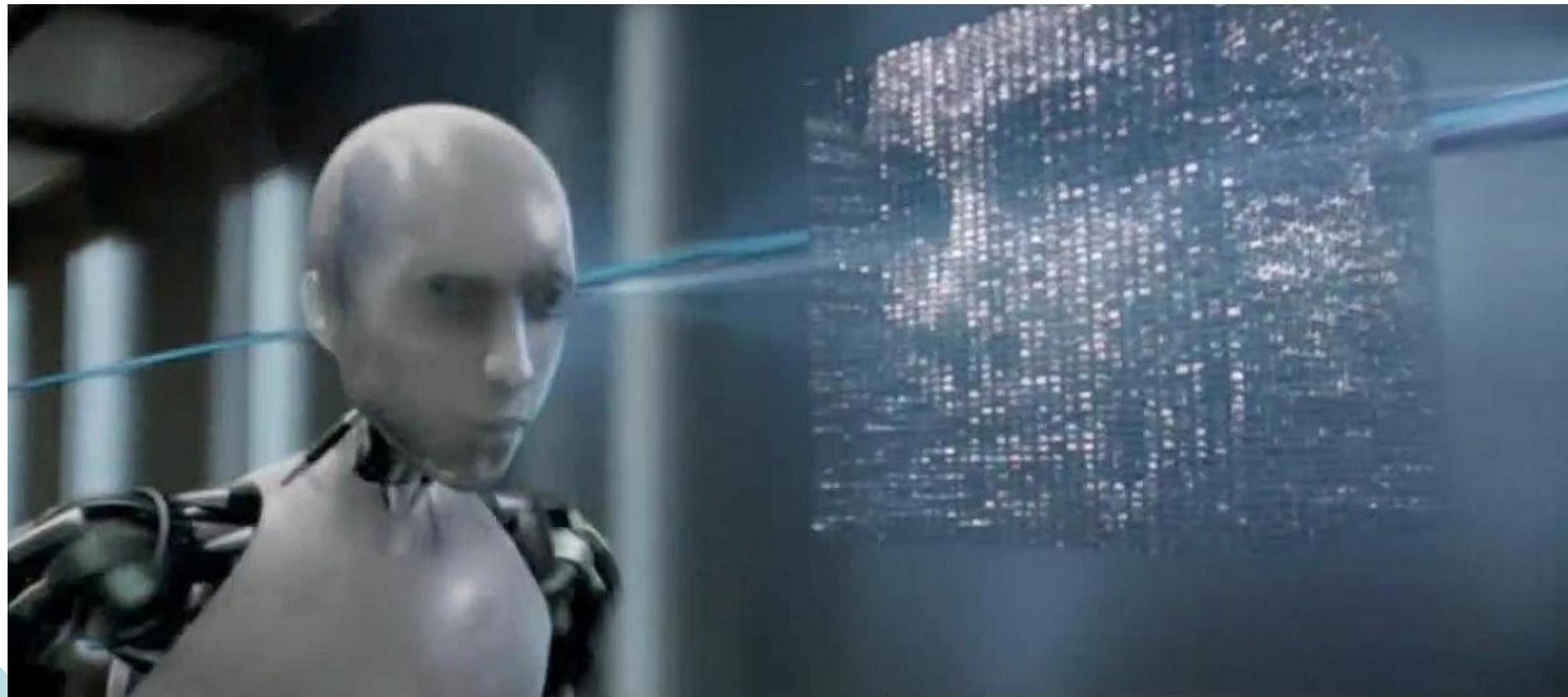
Session Break

AI

What is Artificial Intelligence? How Does AI Work?



AI 란 무엇일까?



자연현상에서 출발

Sidewinder

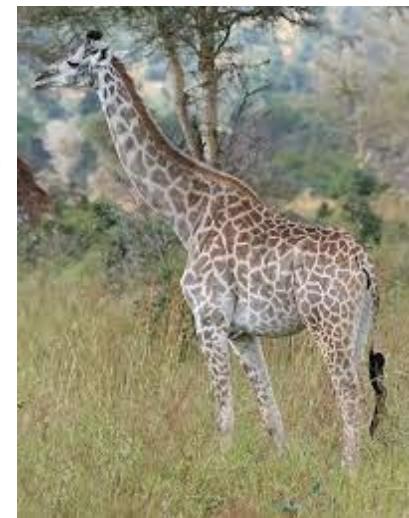


© 2003 JASON JONES

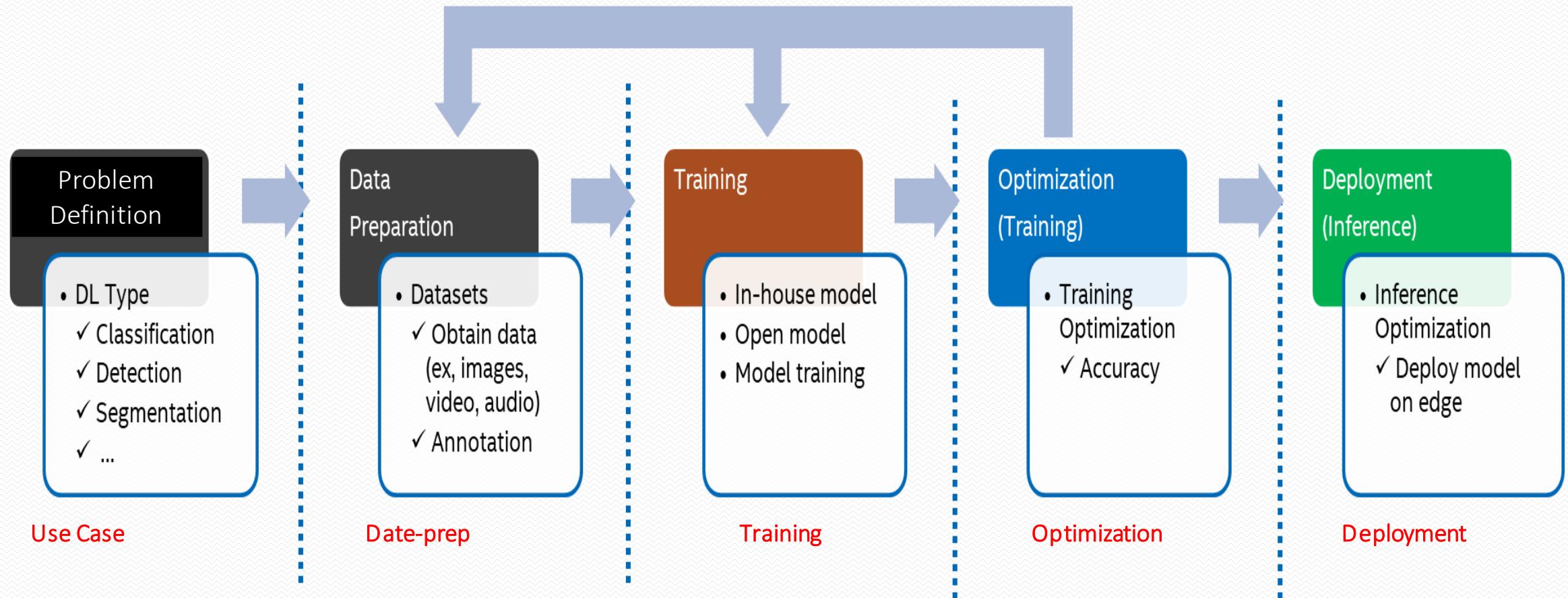
Ultrasound



우리는 어떻게 생각했나



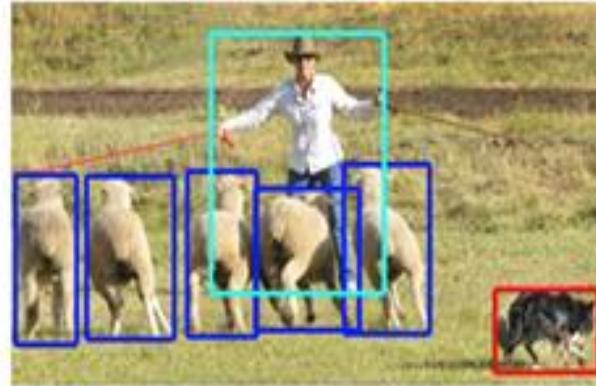
Deep Learning Workflow (생각하게 만들기)



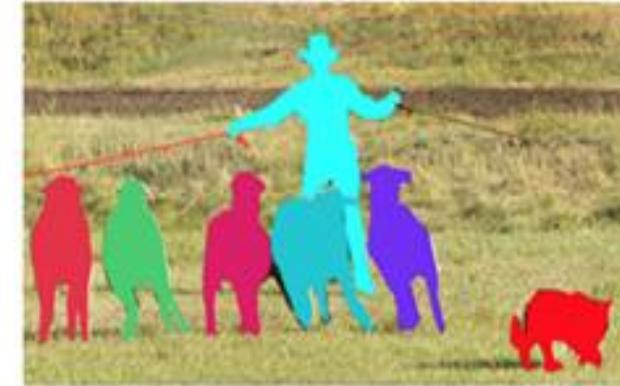
1) Problem Definition



CLASSIFICATION



DETECTION



SEGMENTATION



NEURAL STYLE TRANSFER



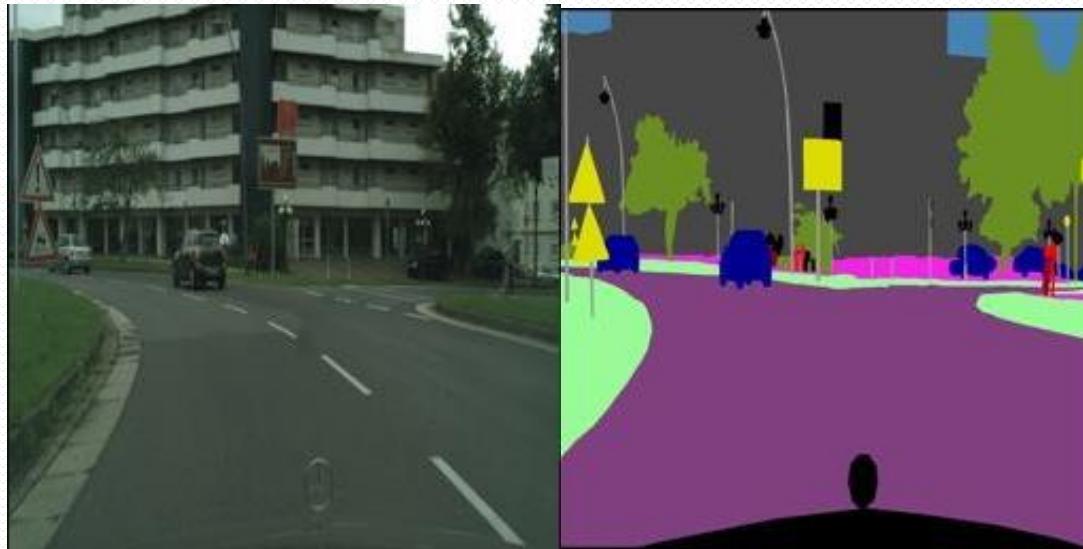
IMAGE CAPTION



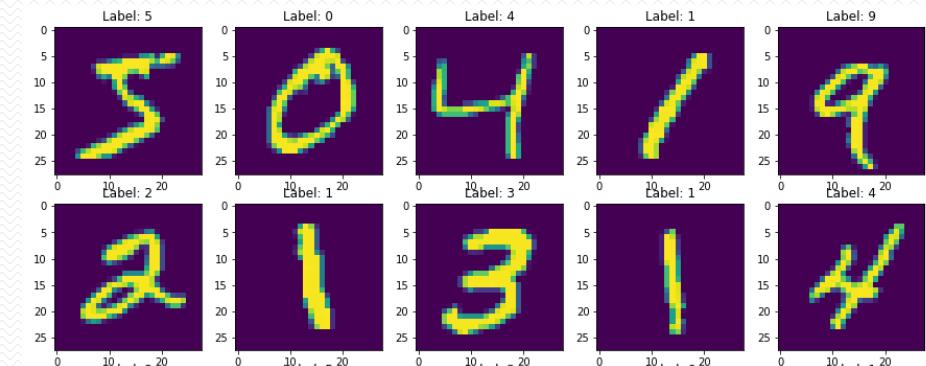
COLORIZATION

2) Dataset preparation/annotation

- Cityscape/Camvid
Semantic Segmentation



- MNIST
Hand-written digit



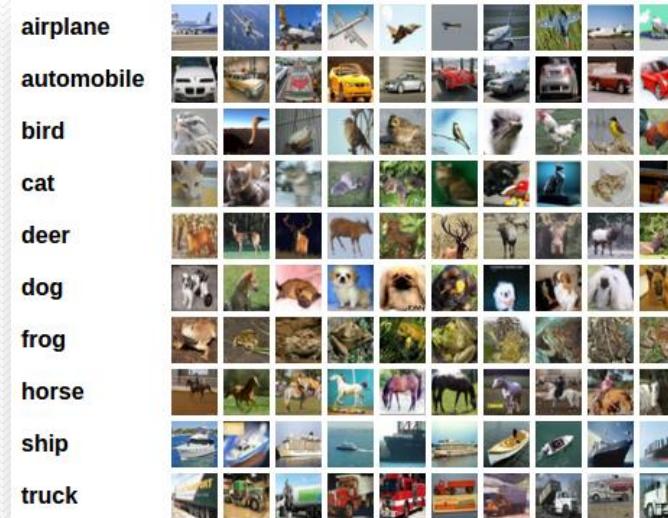
2) Dataset preparation/annotation – cont'd

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

Data & Labels

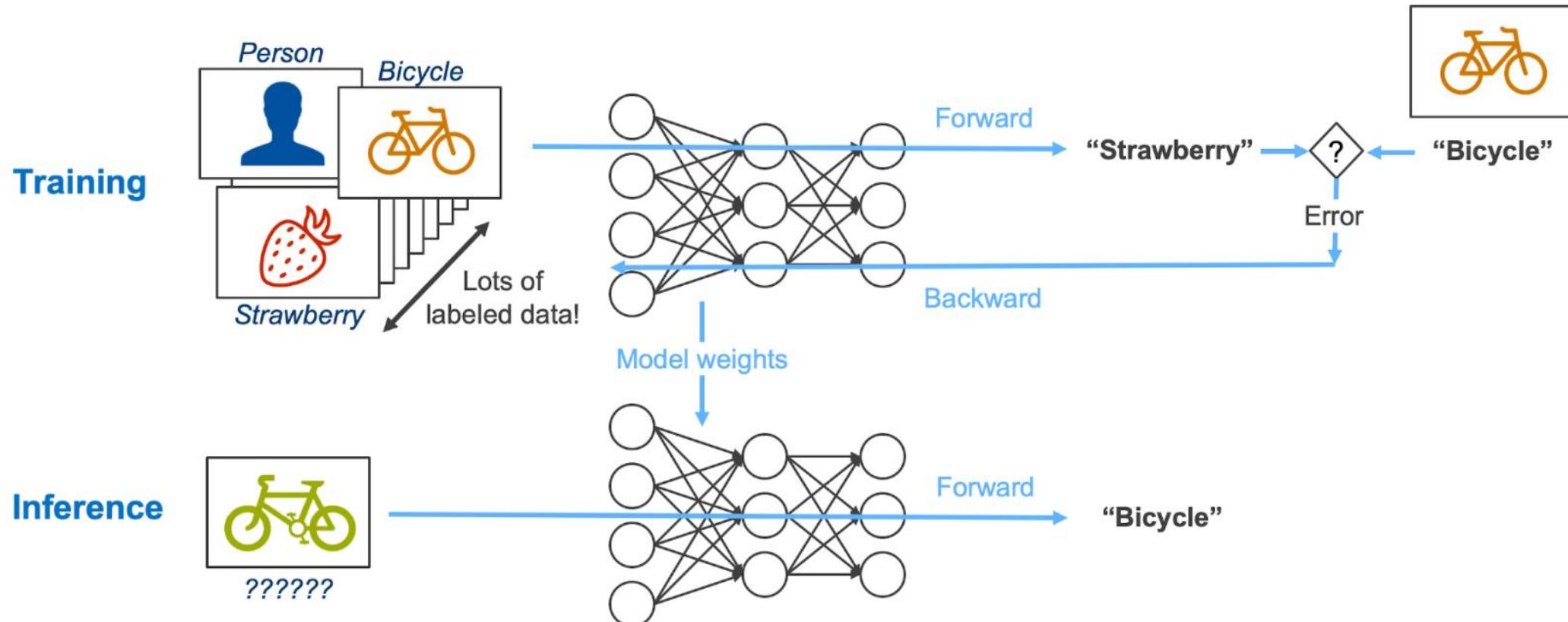
- CIFAR-10 dataset
 - Label : 10 classes
 - Consists of 60,000 (32x32) color images, with 6000 images per class.



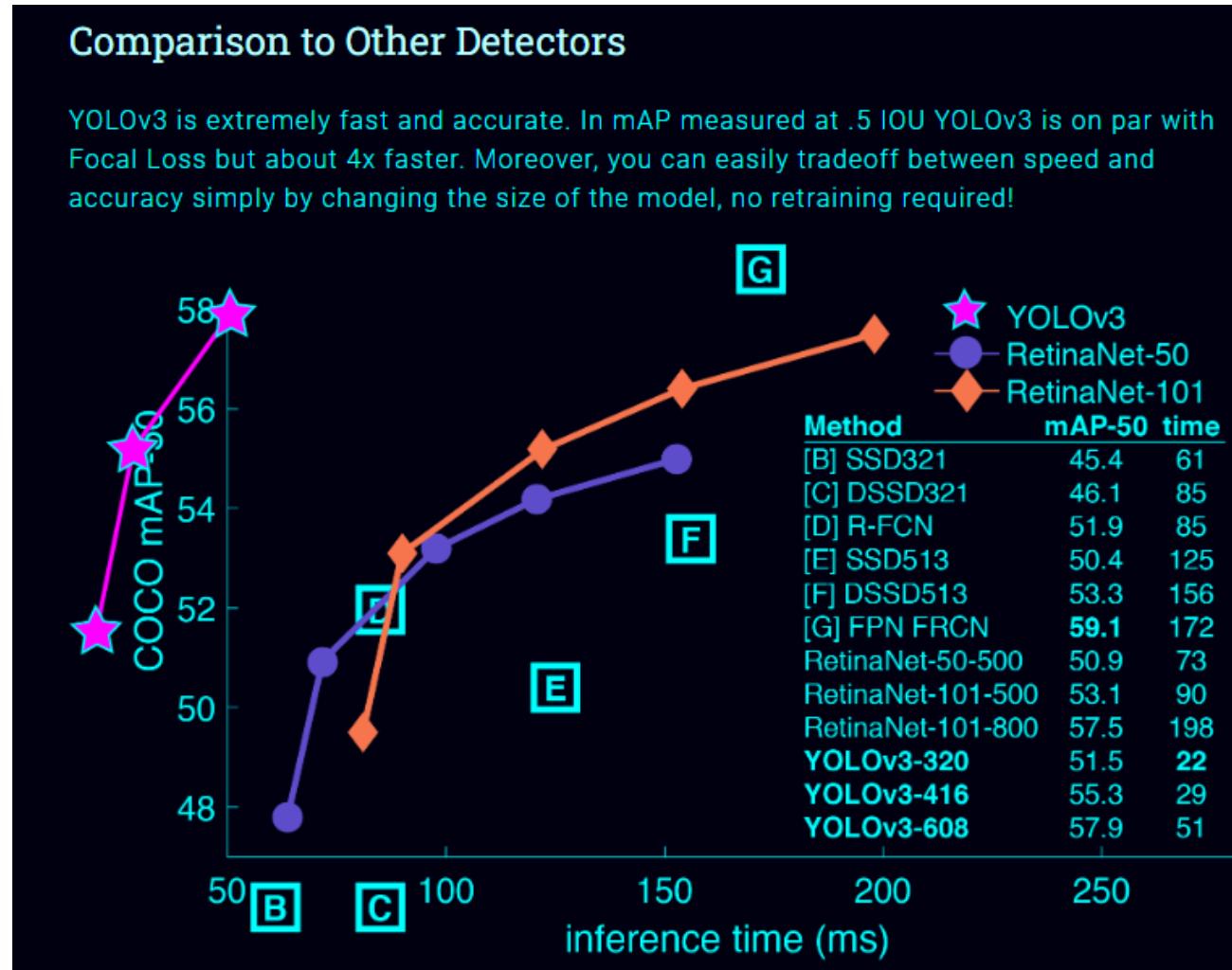
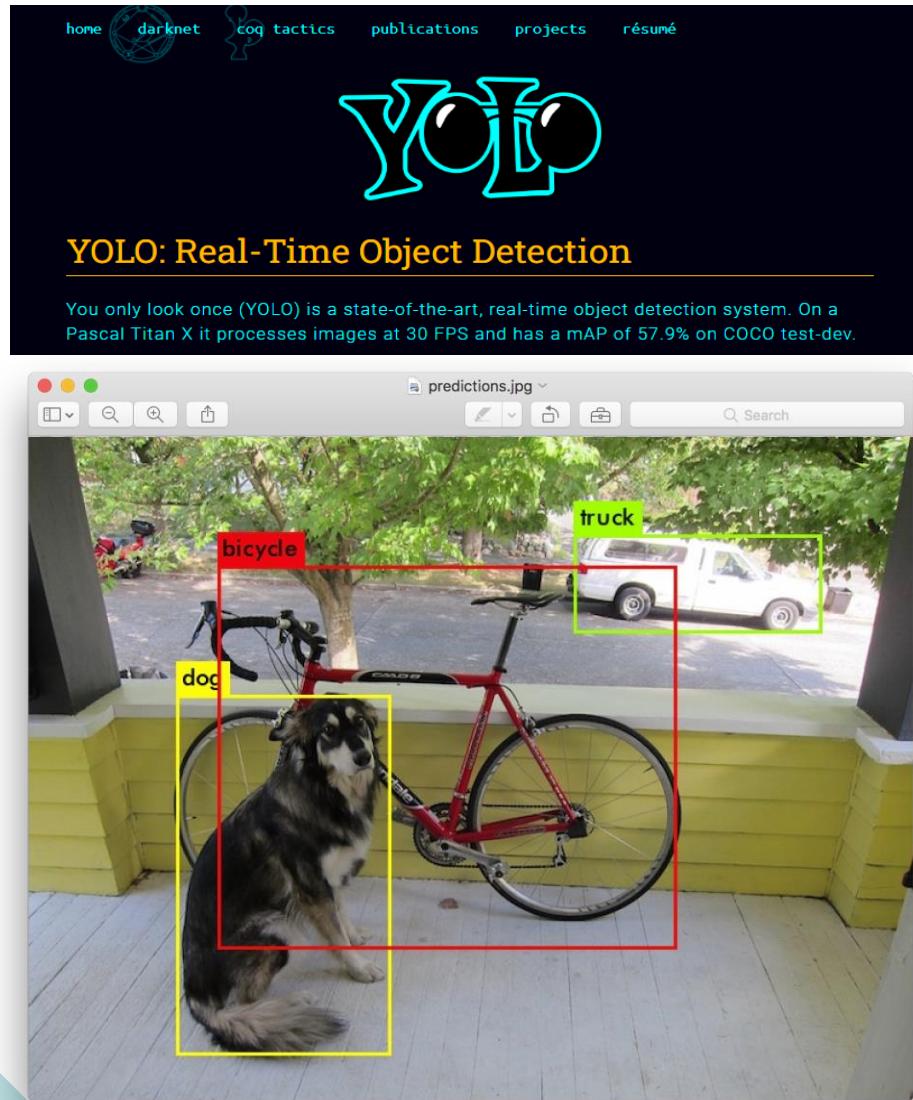
- IMAGENET
 - Label : 1000 classes
 - Consists of 800,000 (256x256 or 224x224) color Images, with 800 images per class.



3) Model selection & Training



Deep Learning Model (YOLO)



4) Optimization

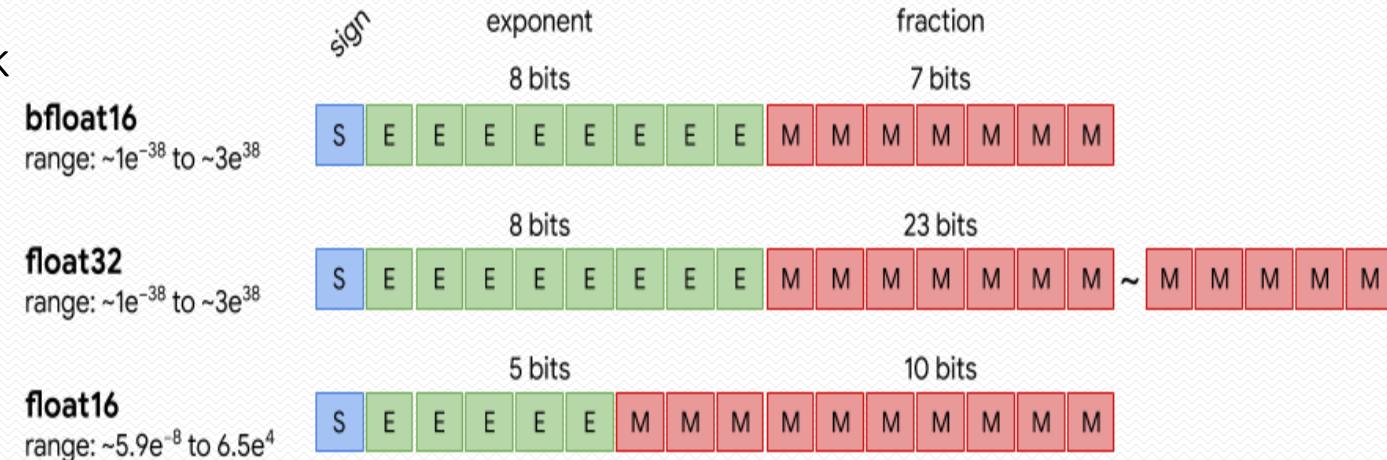
FP16 VS. FP32

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

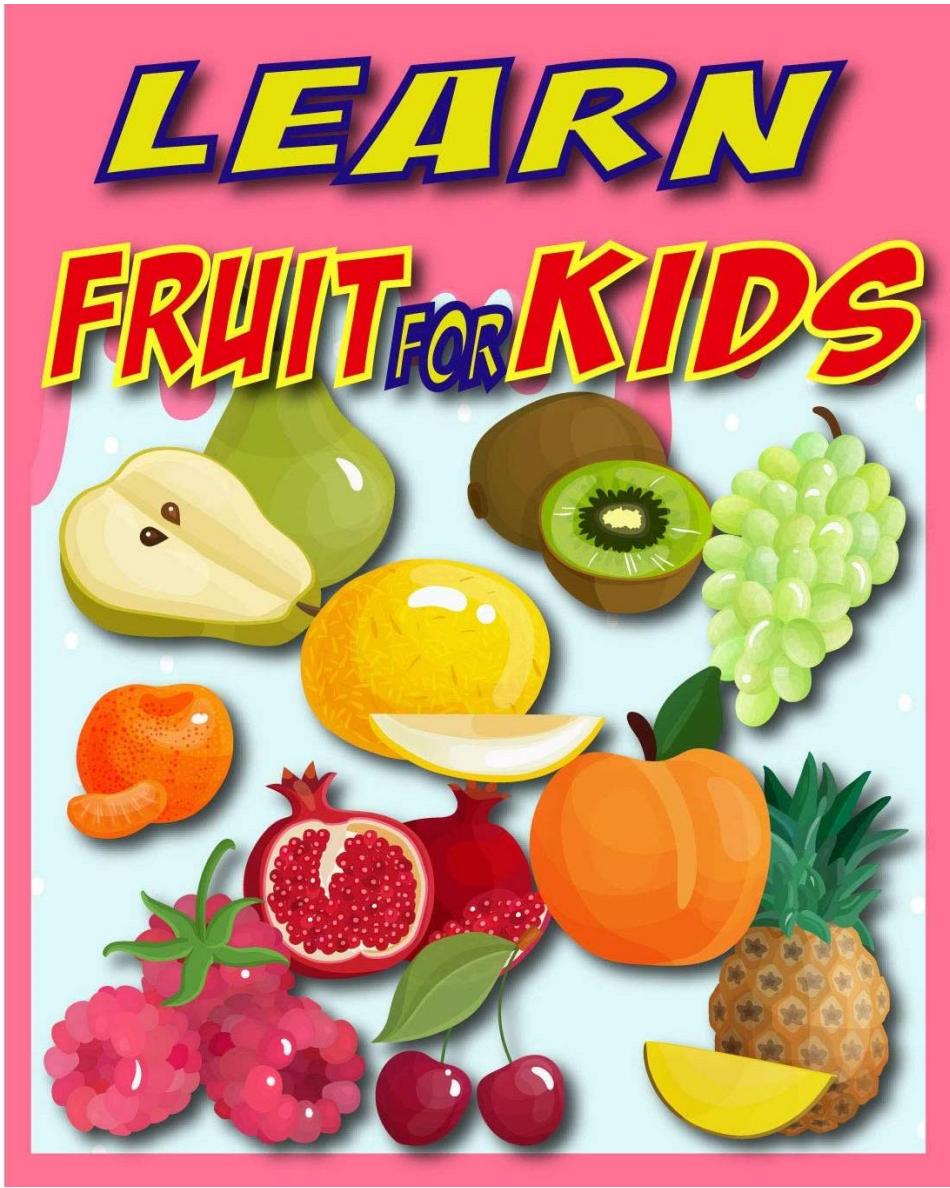
However, this comes at the cost of a small reduction in accuracy



사과!



??



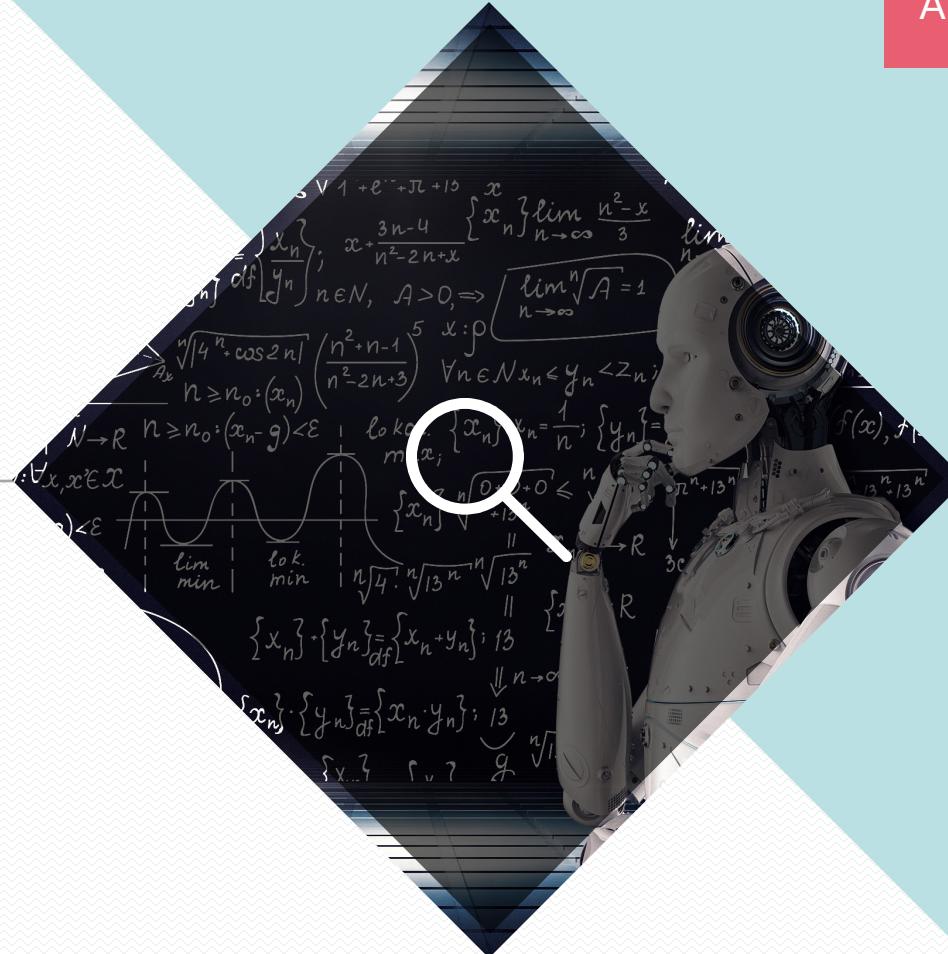
Apple!



Pear!

History

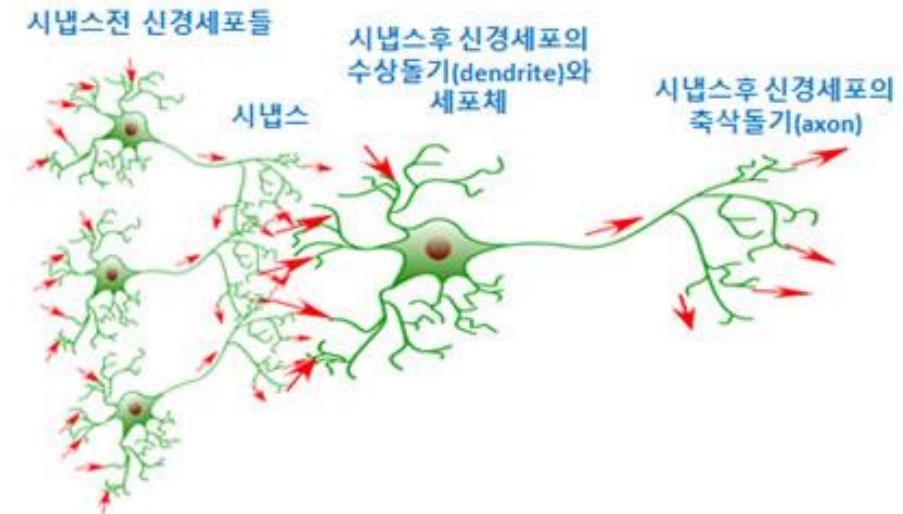
What is Artificial Intelligence? How Does AI Work?



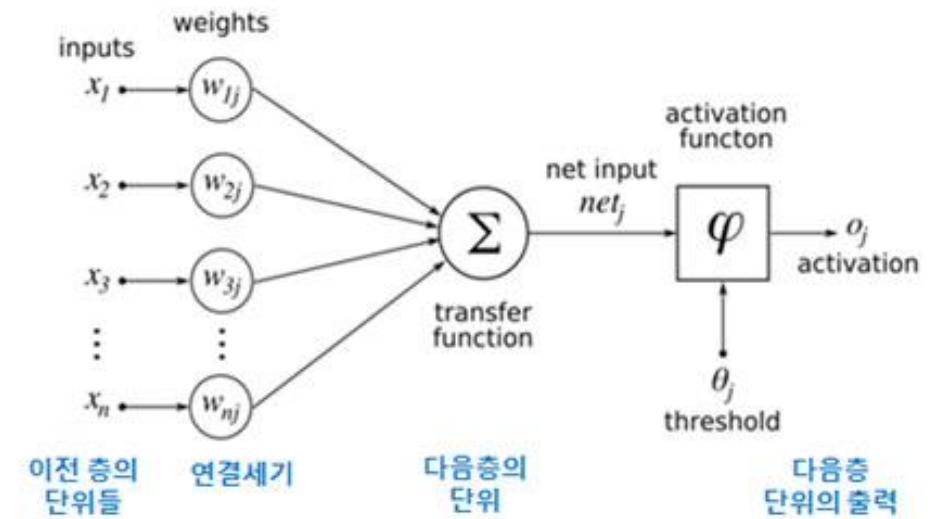
신경망



뇌 속 신경망

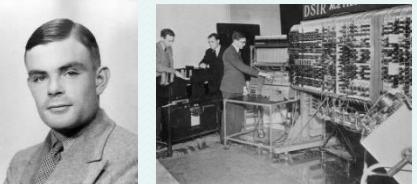


인공 신경망



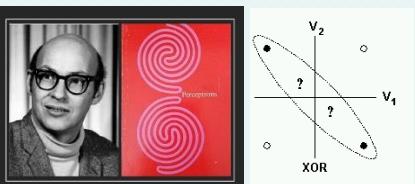
HISTORY OF A.I.

Alan Turing
Turing Test



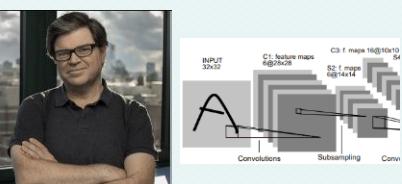
1950

Marvin Minsky
A book of "Perceptrons"



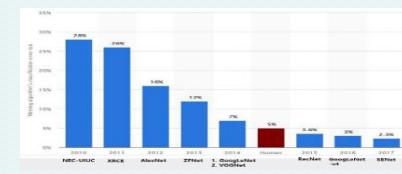
1969

Yann Lecun
LeNet-5



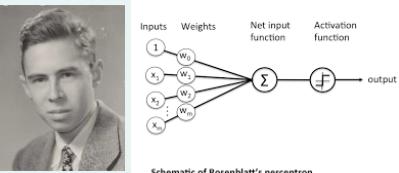
1998

ResNet
Accuracy is higher than human



2015

Frank Rosenblatt
Perceptron

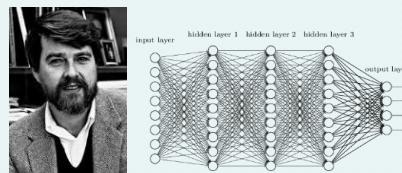


1958

AI Winter

1986

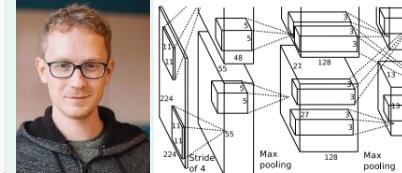
David Rumelhart
Multi Layer Perceptron (MLP)



AI Winter

2012

Alex Krizhevsky
AlexNet

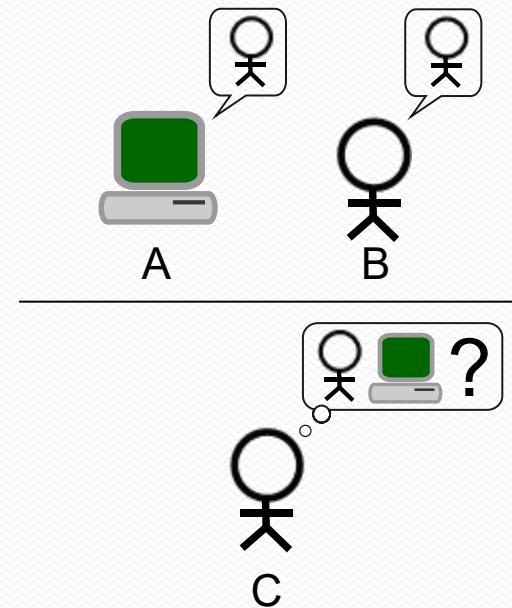


2016

Google Brain
AlphaGo



HISTORY OF A.I.

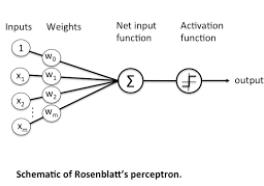


- 튜링 테스트(Turing test)의 의미
 - 기계가 생각하는 것이 가능한가?
 - Computing Machinery and Intelligence 논문
- 각각 분리된 채, 하나는 인간, 다른 하나는 기계와 대화를 하는 Test
- 5분동안 기계와 대화한다는 것을 인지하지 못하면 성공
- 대화에서 질문에 올바른 대답을 하는것 보단 얼마나 인간같이 대화할 수 있느냐로 판단

HISTORY OF A.I.

1958

Frank Rosenblatt
Perceptron



- 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*



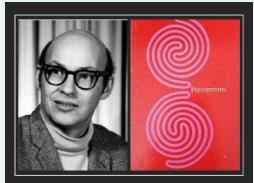
If we are eventually to understand the capability of higher organisms for perceptual recognition, generalization, recall, and thinking, we must first have answers to three fundamental question.

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?*

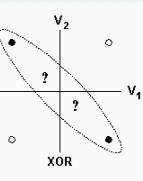
HISTORY OF A.I.

Marvin Minsky

A book of "Perceptrons"

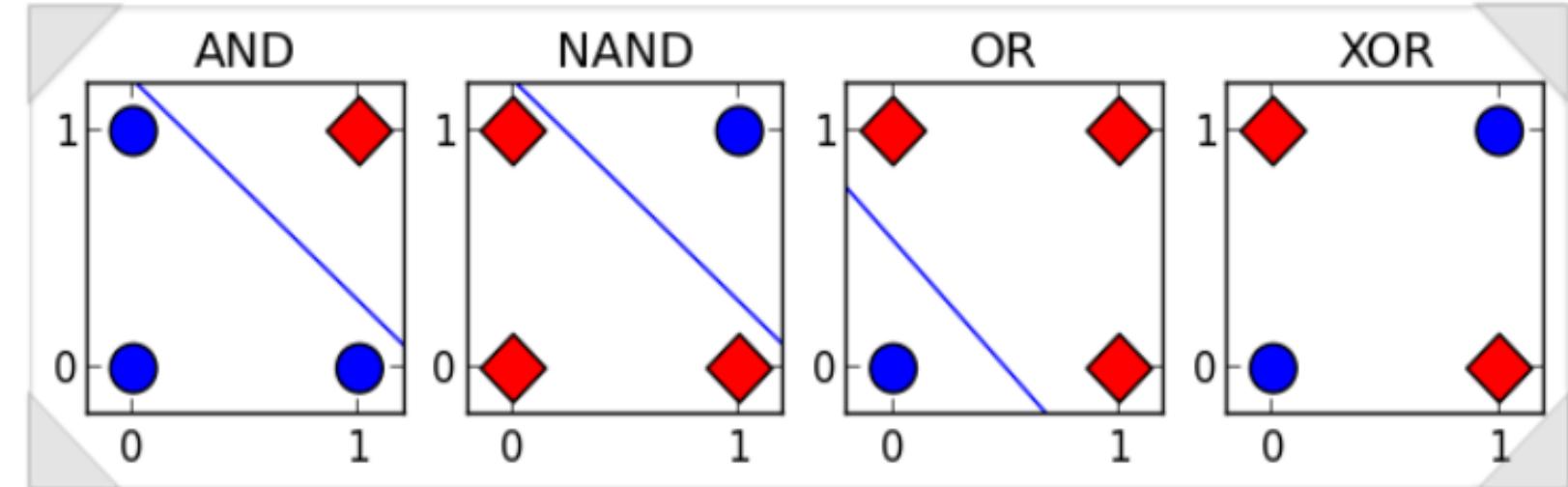


1969

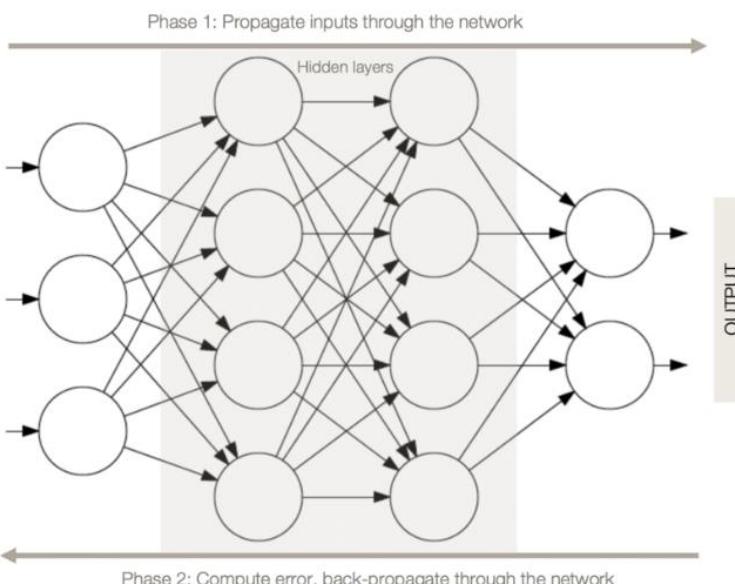
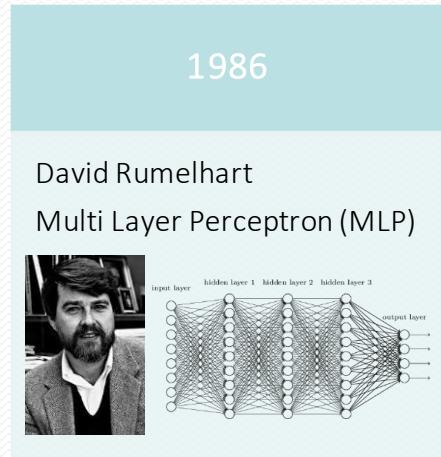


마빈 민스키

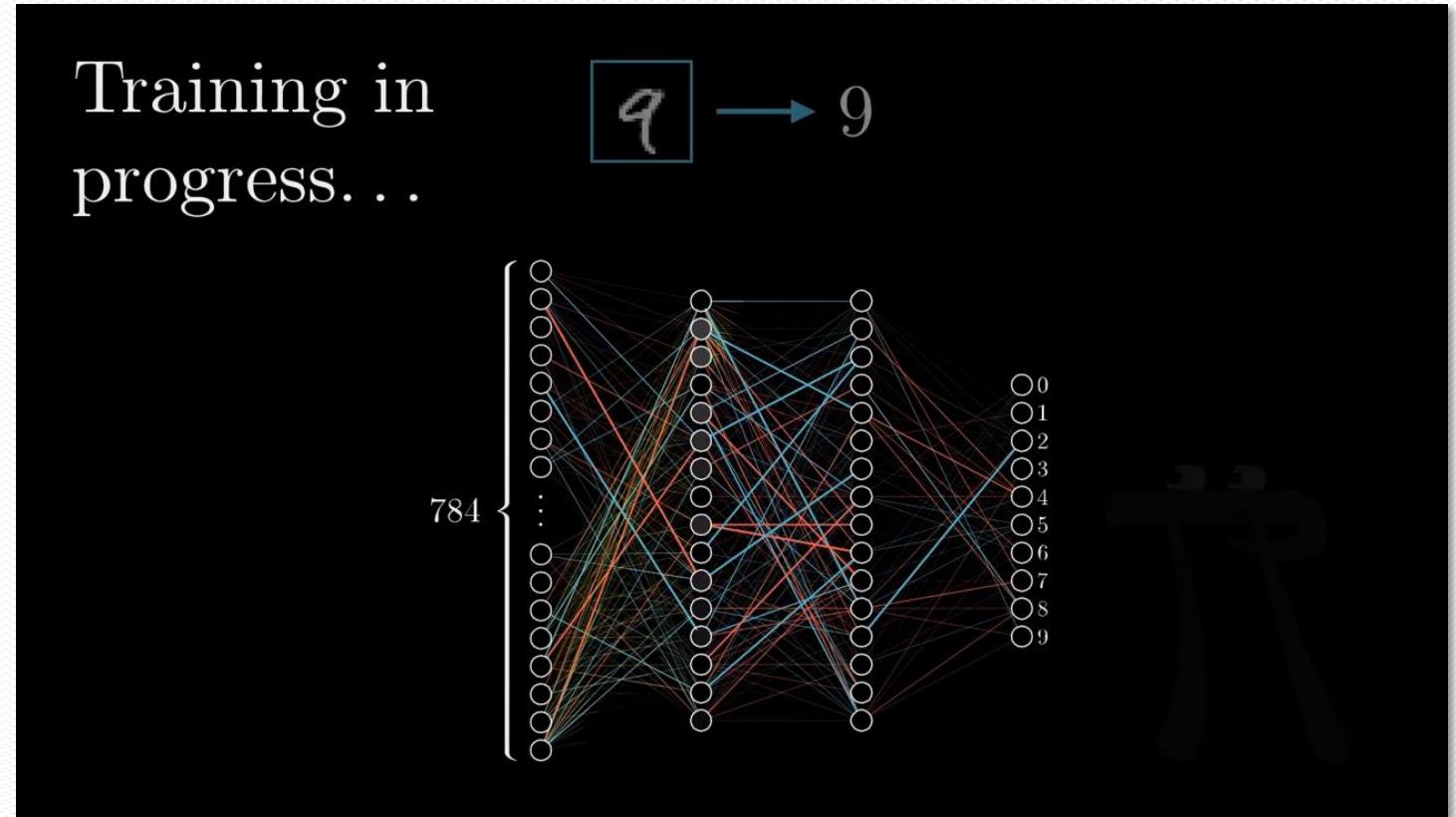
- *Perceptrons* – Minsky & Papert 출간, 퍼셉트론 모델에 대해 철저한 분석 및 그 한계에 대해서도 논리정연하게 분석



HISTORY OF A.I.



- 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.

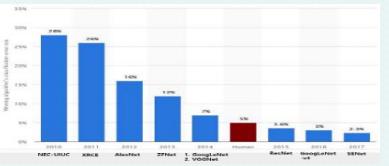
Yann Lecun
LeNet-5



INPUT 32x32
E1: 16 maps 8x28x28
C3: 16@10x10 S4
E2: 16 maps 8x14x14
S2: 16@14x14
Convs
Subsampling
Conn

1998

ResNet
Accuracy is higher than human

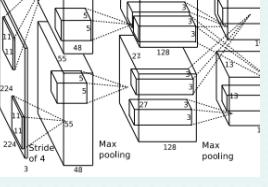


Model	Accuracy (%)
Net-5	24.0%
LeNet	24.0%
Baseline	19.0%
2011	17.0%
2012	13.0%
2013	11.0%
2014	10.0%
2015	9.0%
ResNet	15.0%
Baseline	9.0%
2016	9.0%
2017	9.0%
2018	9.0%
2019	9.0%
2020	9.0%
Human	9.0%

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

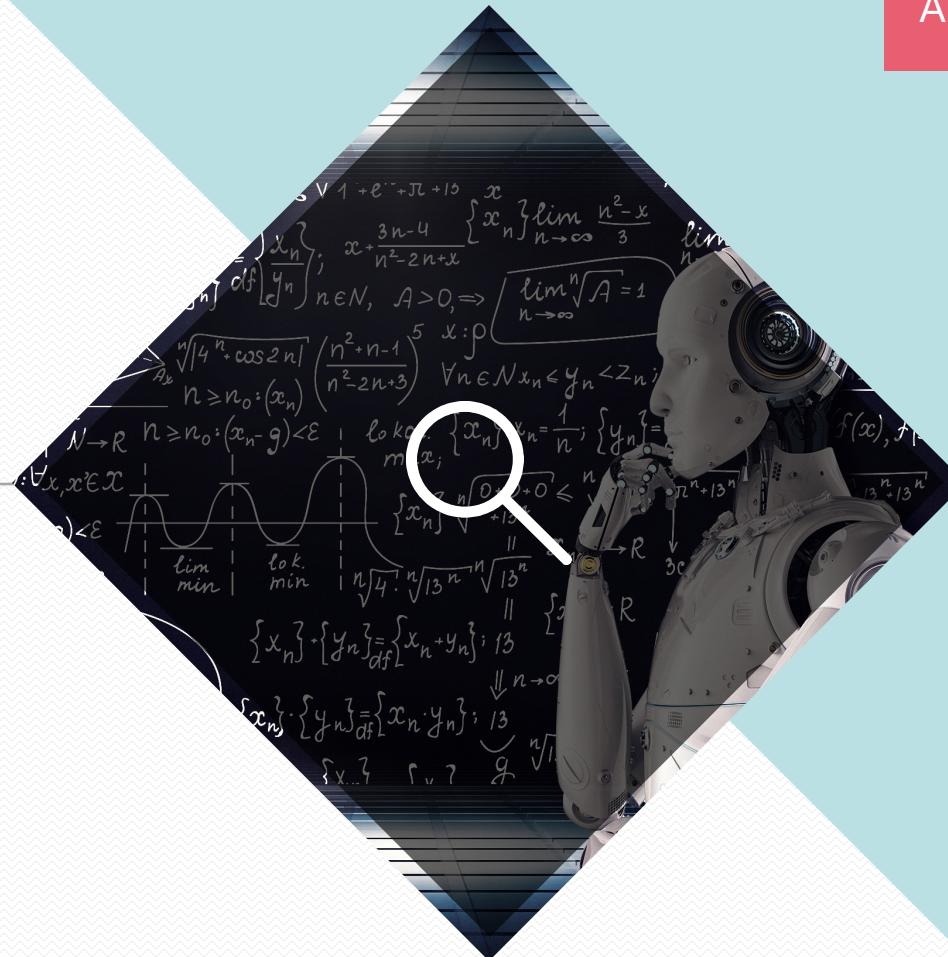
ResNet

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

Session Break

ANN

ARTIFICIAL NEURAL NETWORKS



CONTENTS

NEURAL NETWORKS

- WHAT IS THIS?
- NEURON / WEIGHTS / BIASES

PERCEPTRON

- LOGICAL GATE EXAMPLE
- MULTI-LAYER PERCEPTRON , MLP

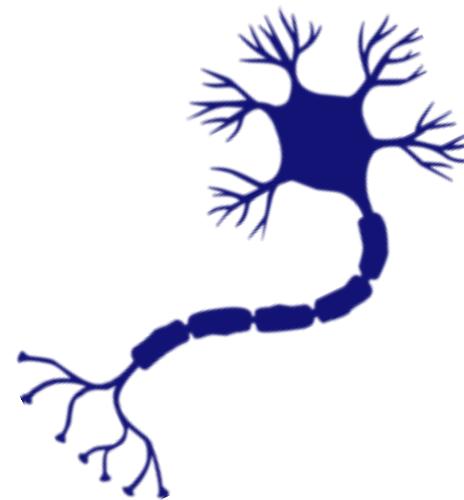
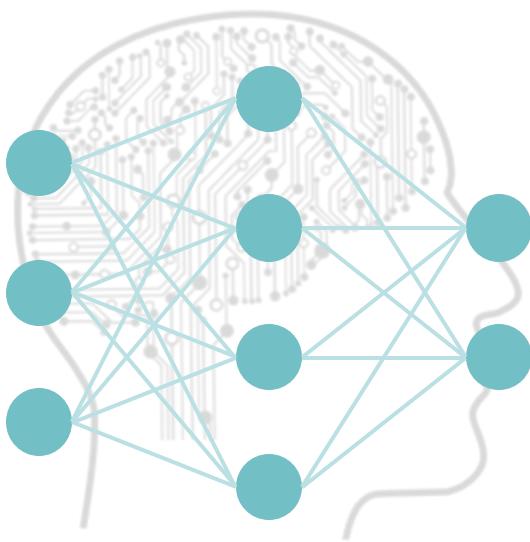
TECHNICAL DEEP DIVE

- ACTIVATION FUNCTIONS
- Optimizer

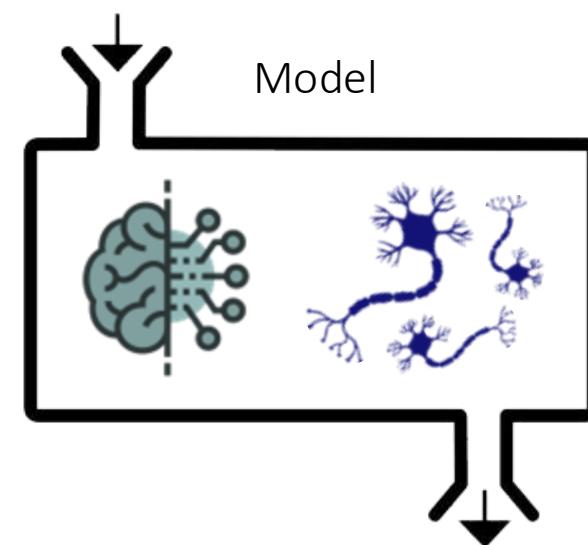
Hands-on

NEWRAL NETWORKS OVERVIEW

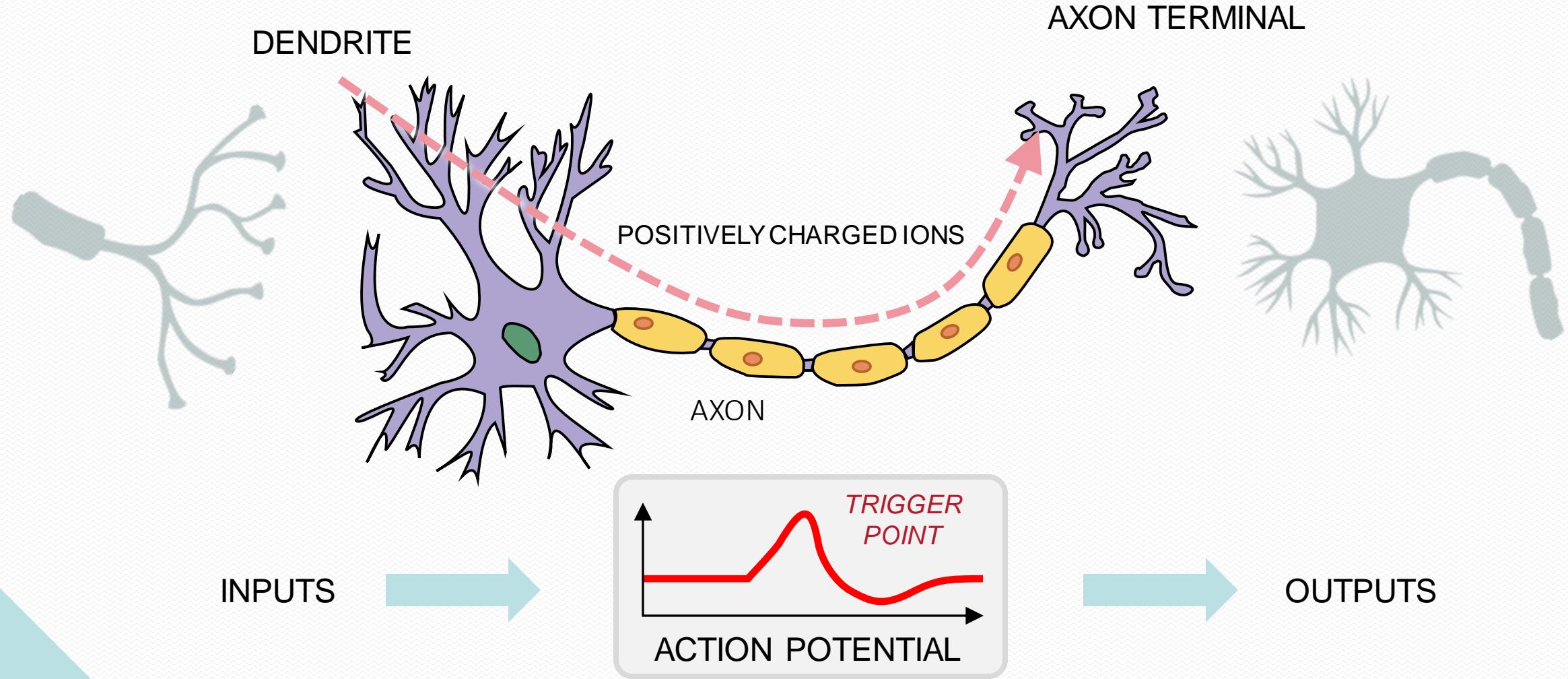
Neural networks



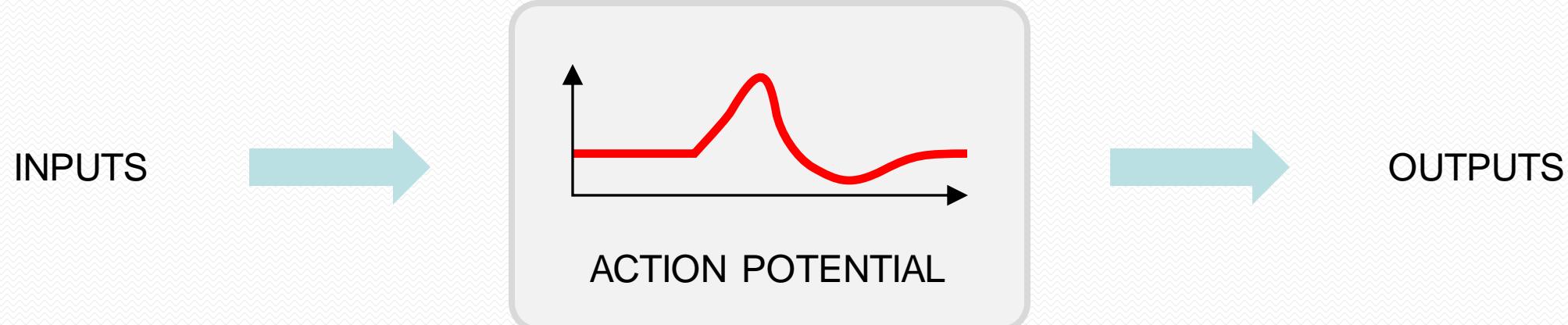
Neuron in human brain



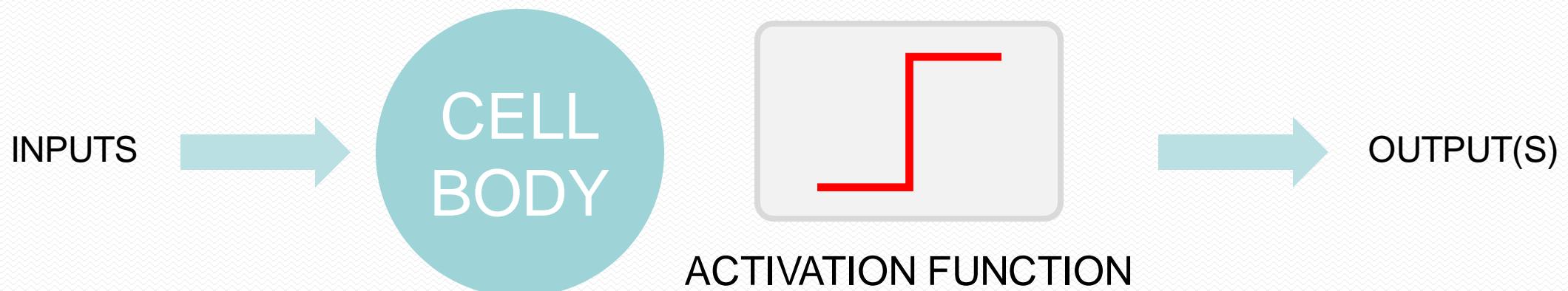
NEURON IN OUR BRAINS



ACTIVATION



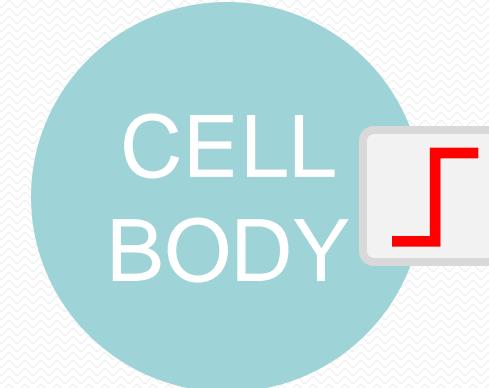
PERCEPTRON



WEIGHTS AND BIASES

WEIGHTS

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



Activation
Function

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

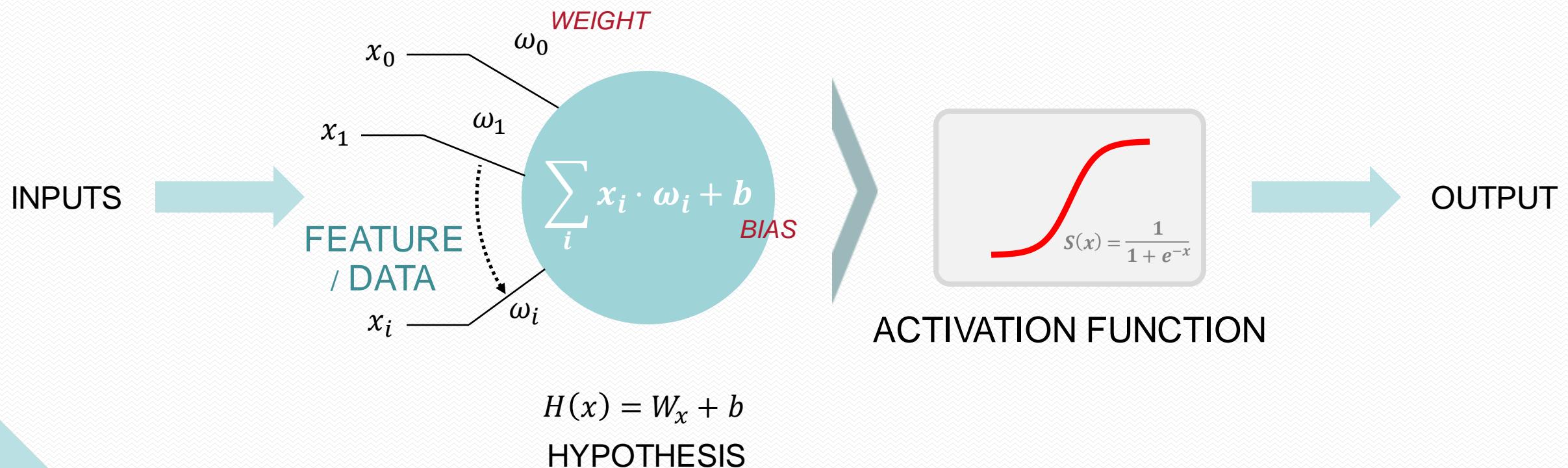
BIASES

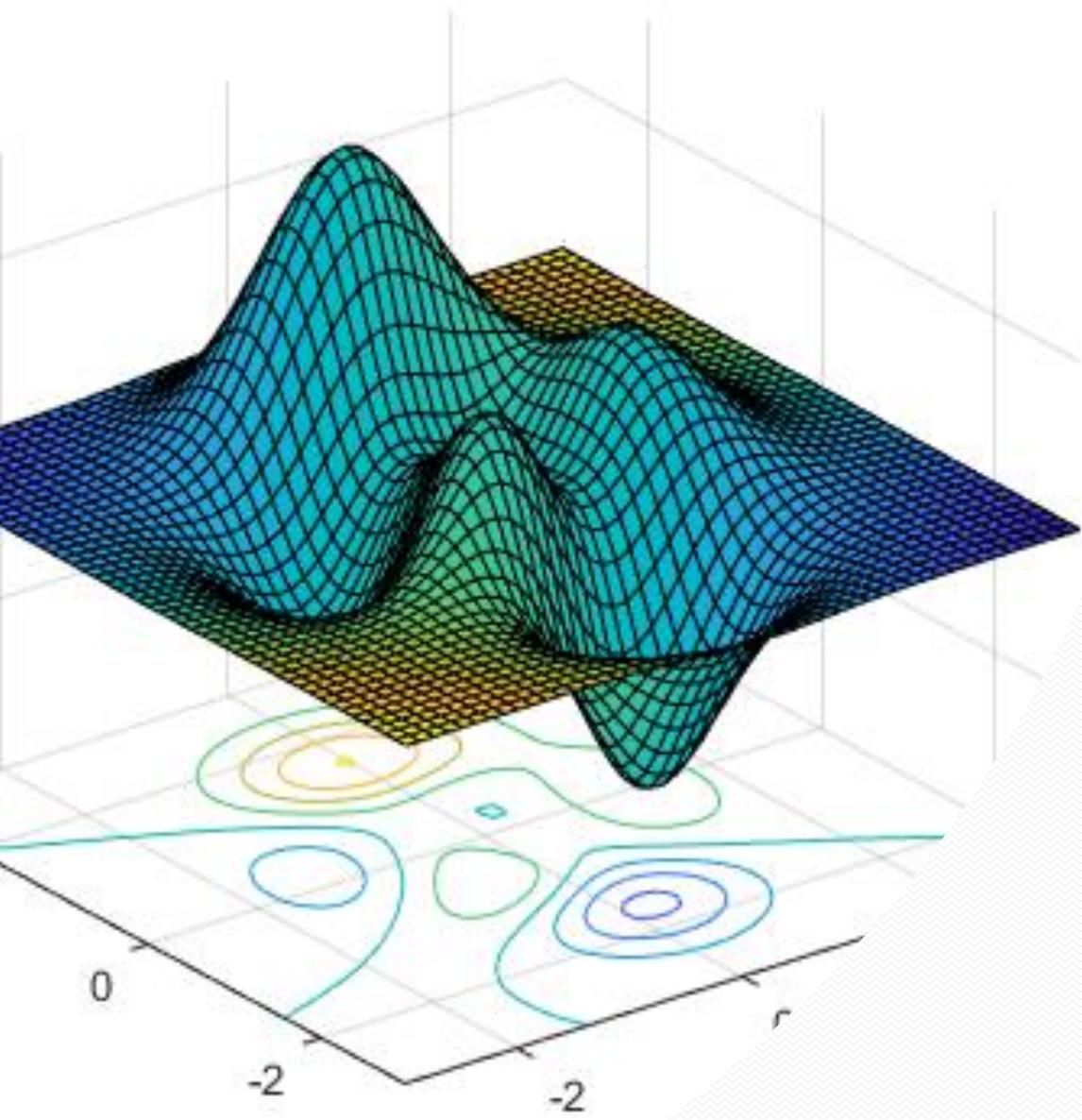
- ADDITIONAL INPUT CONSTANT VALUES FOR NEXT NEURONS
- NOT INFLUENCED BY THE PREVIOUS NEURONS

These are the connections that go out of their own weights

PERCEPTRON

ALGORITHM FOR SUPERVISED LEARNING OF BINARY CLASSIFIERS

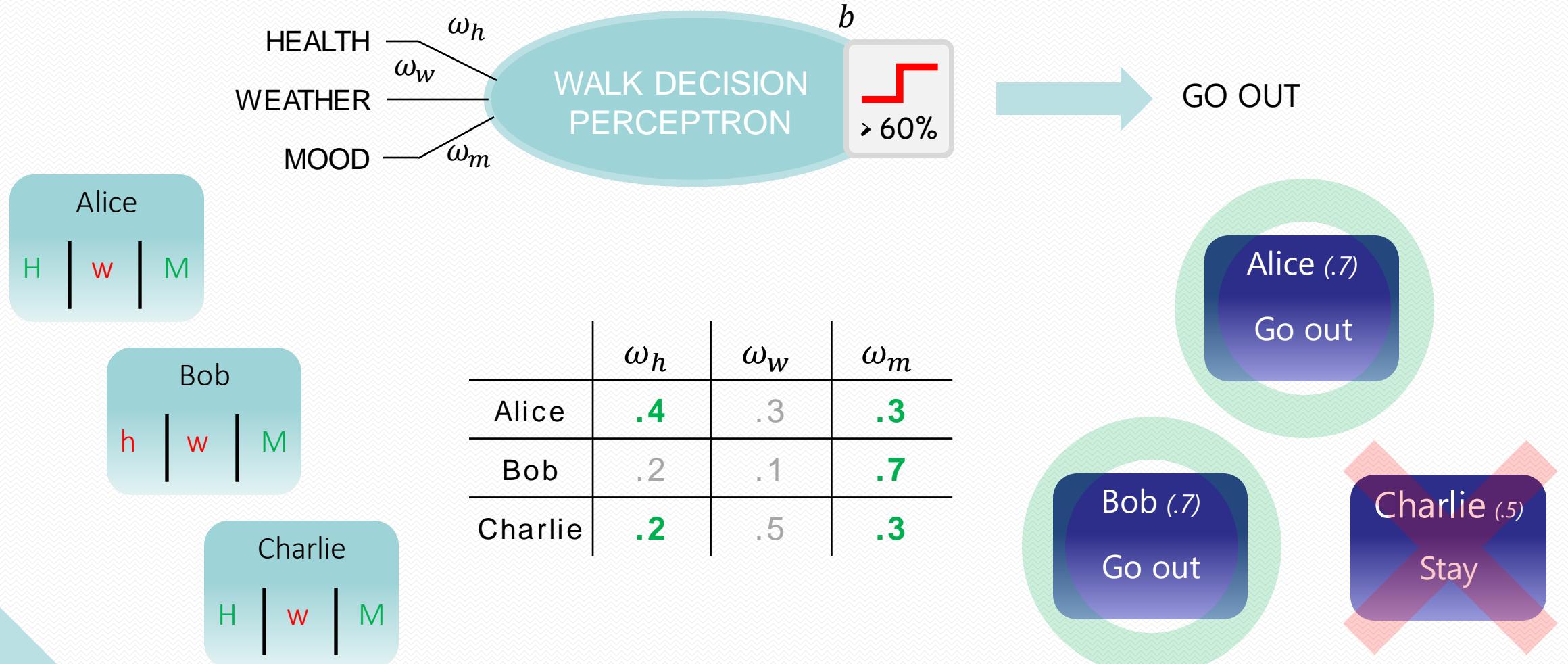




Perceptron

Logic gate

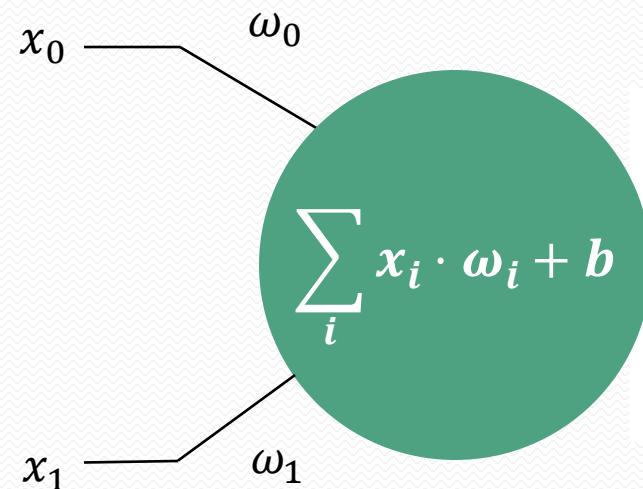
PERCEPTRON EXAMPLE



Perceptron Example – Logical OR

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

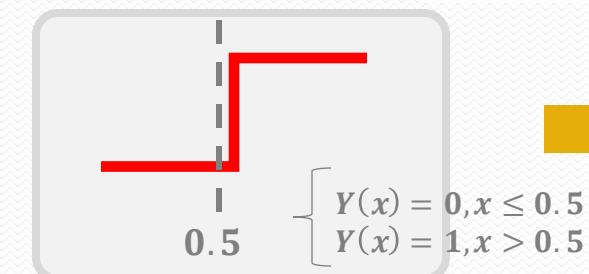
x_0	x_1	w_0	w_1	b	$H(x)$	Output
0	0	0.2	0.2	0.4	0.4	0
0	1				0.6	1
1	0				0.6	1
1	1				0.8	1



⋮

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



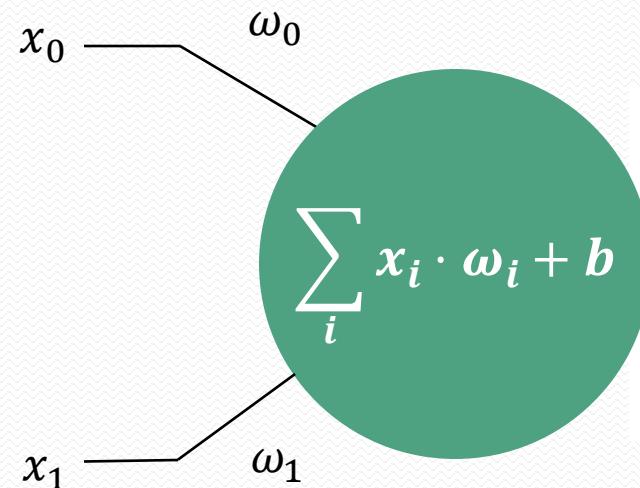
Activation Function



Perceptron Example – Logical NAND

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

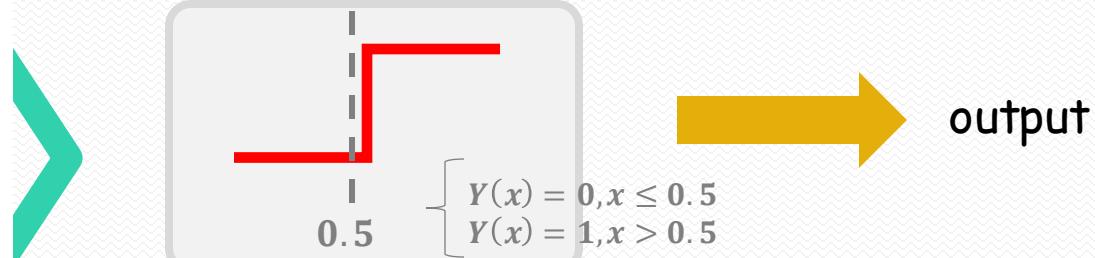
x_0	x_1	w_0	w_1	b	$H(x)$	Output
0	0	-0.5	-0.5	1.1	1.1	1
0	1				0.6	1
1	0				0.6	1
1	1				0.1	0



⋮

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

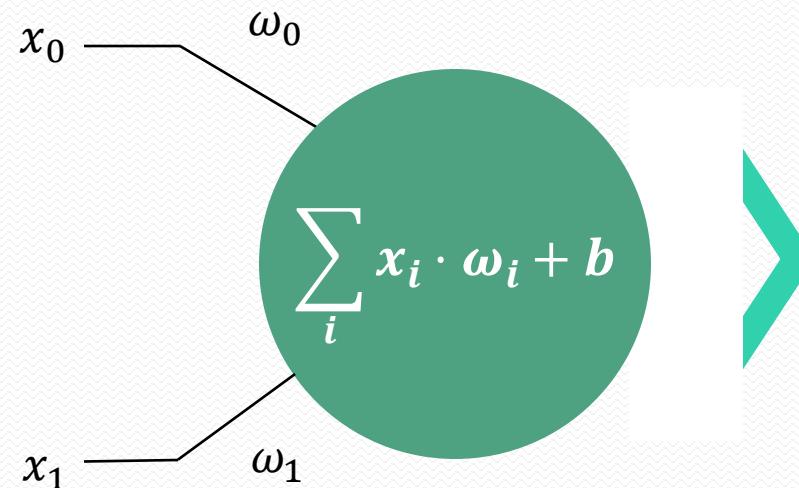


Activation Function

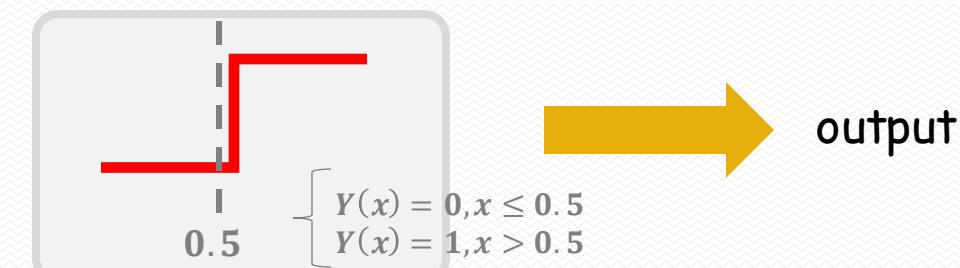
Perceptron Example – Logical XOR

- 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



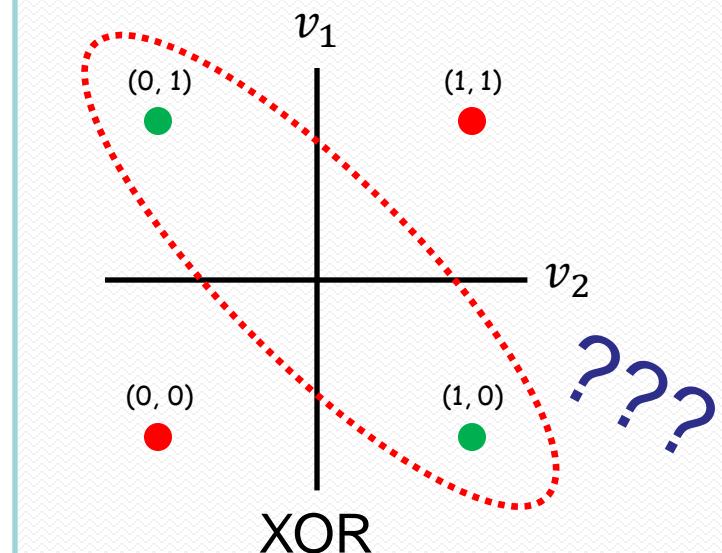
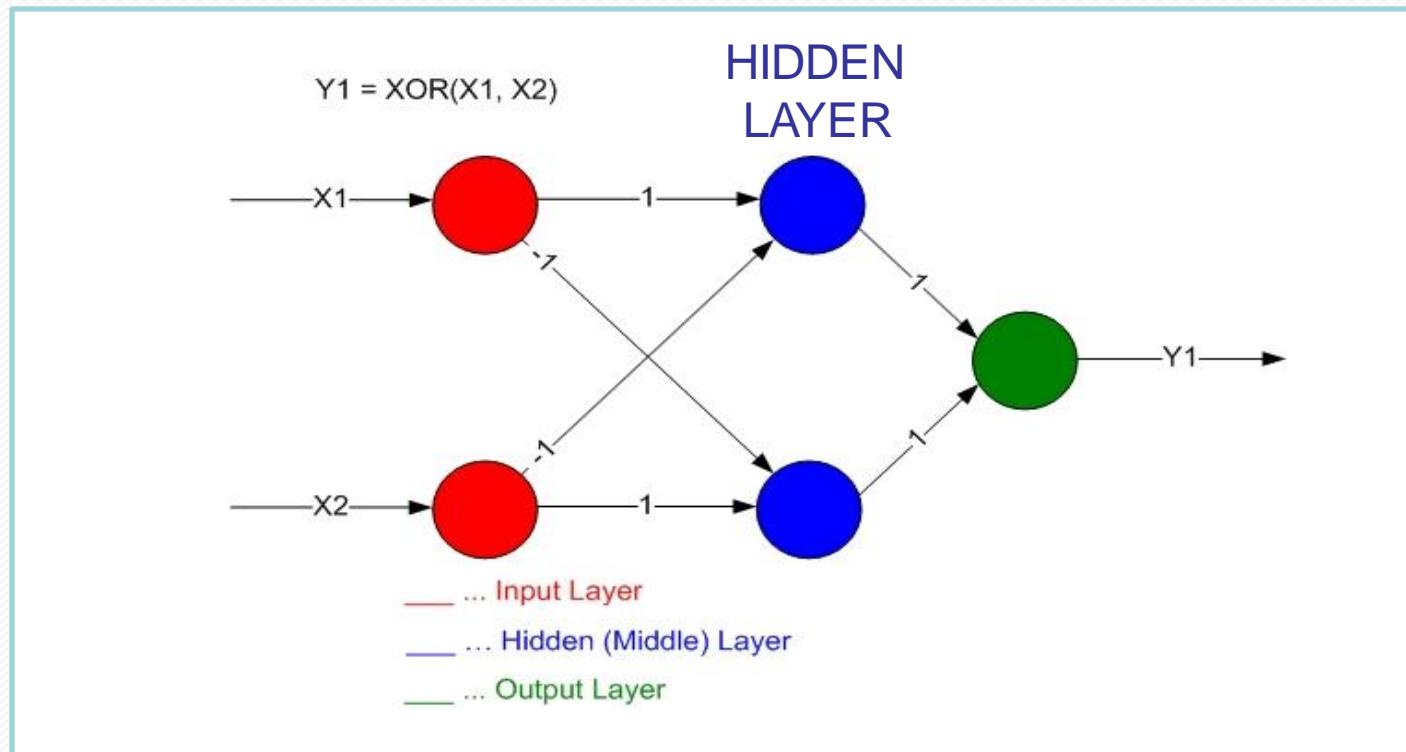
Activation Function

PROBLEM OF PERCEPTRON

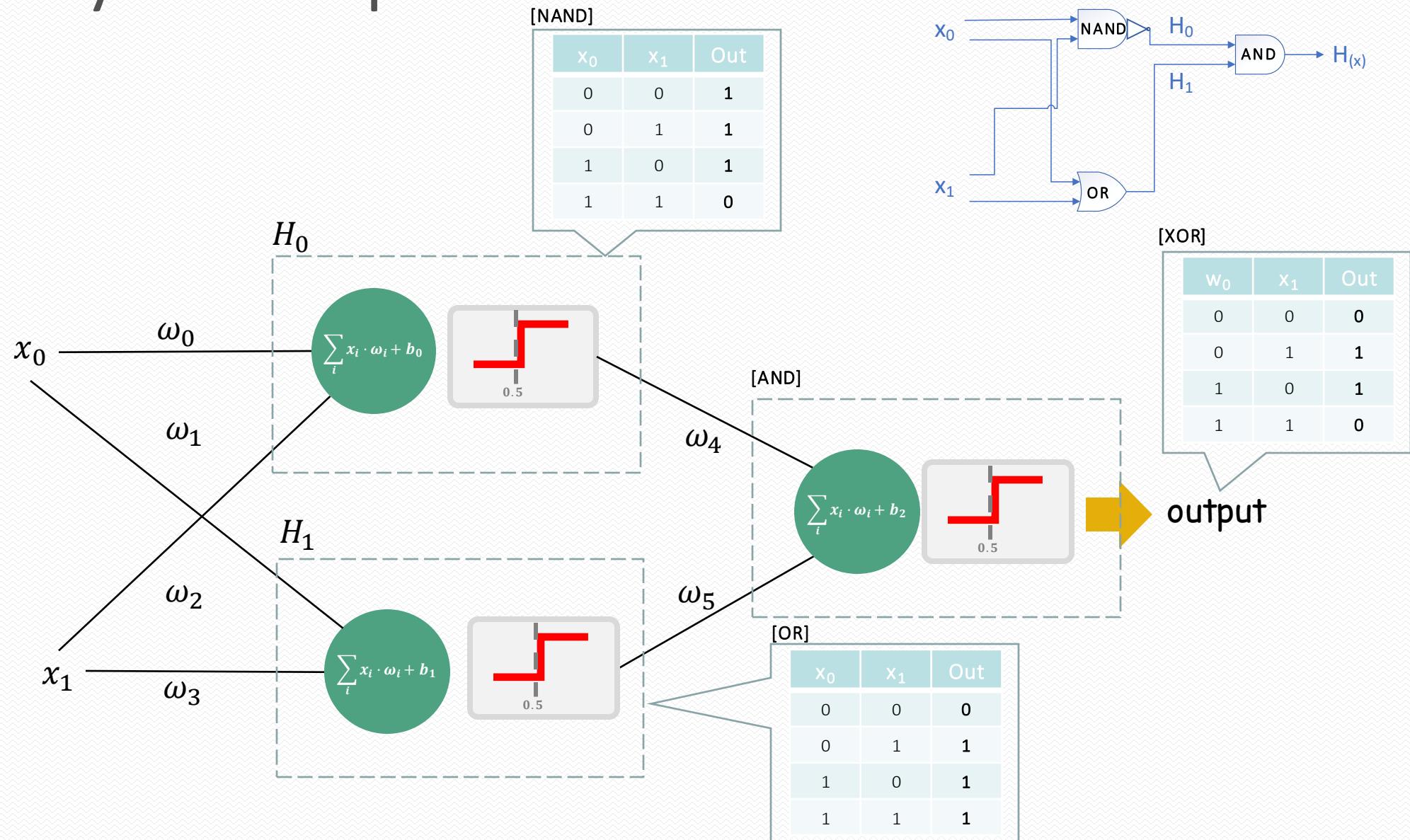
PERCEPTRON IS NOT AN ALL-ROUND PLAYER

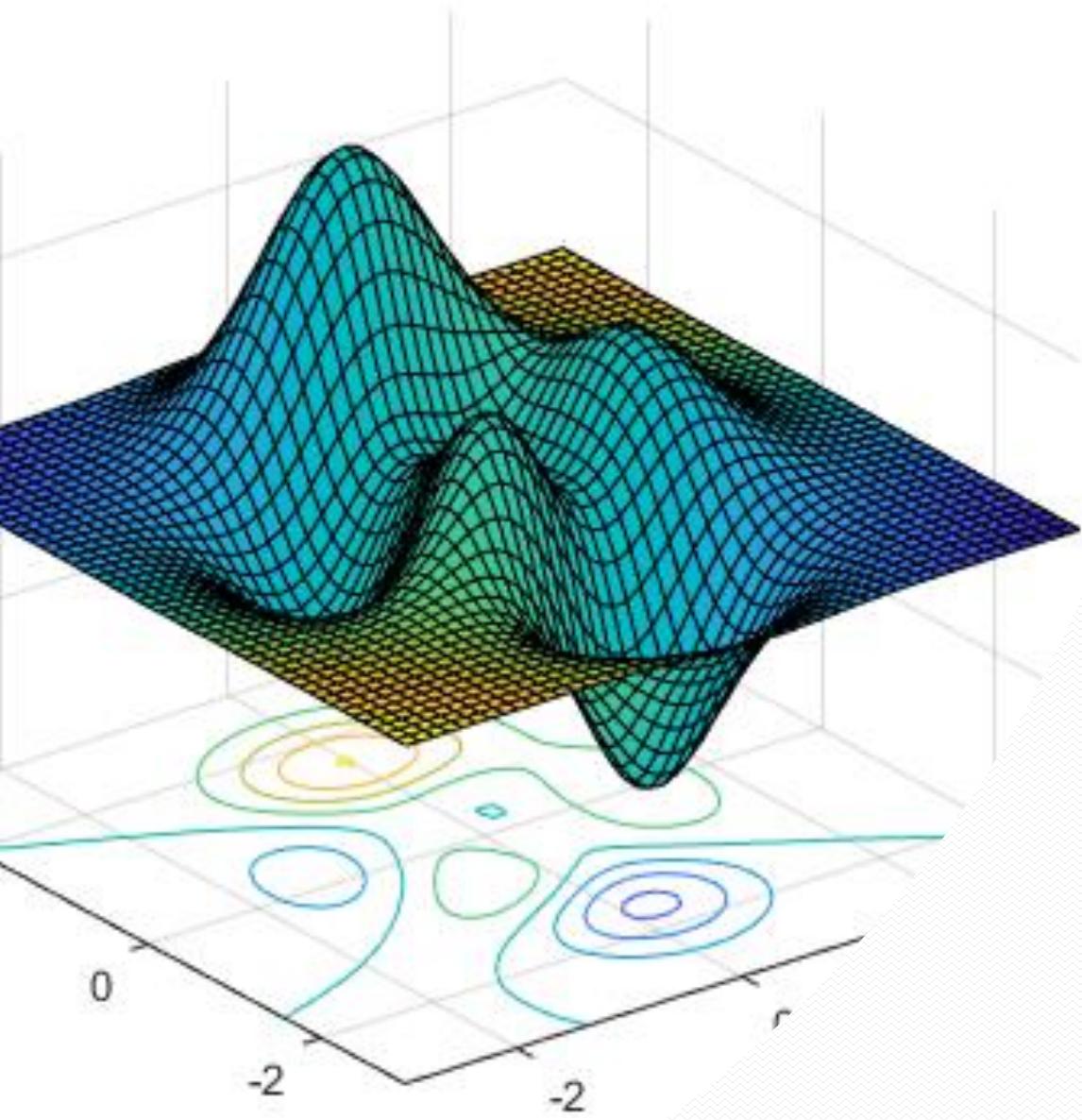
PROBLEMS THAT ARE NOT CLASSIFIED AS LINEAR CANNOT BE CALCULATED

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



Multi Layer Perceptron - XOR



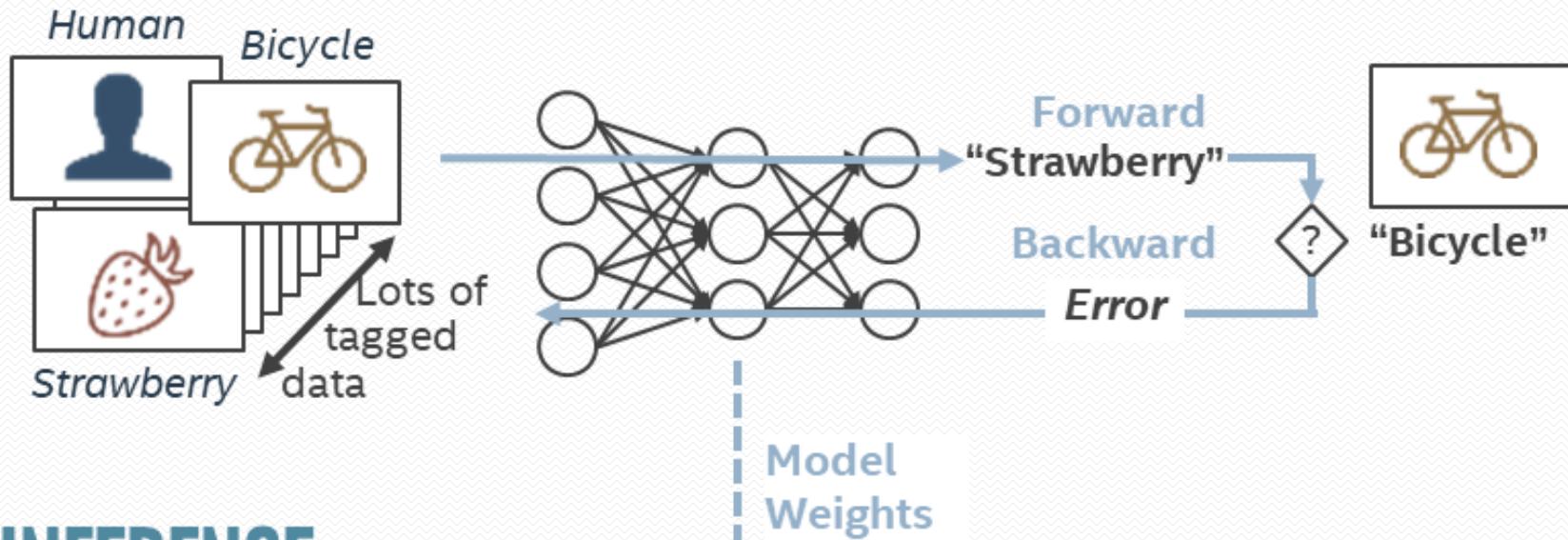


Deep-dive

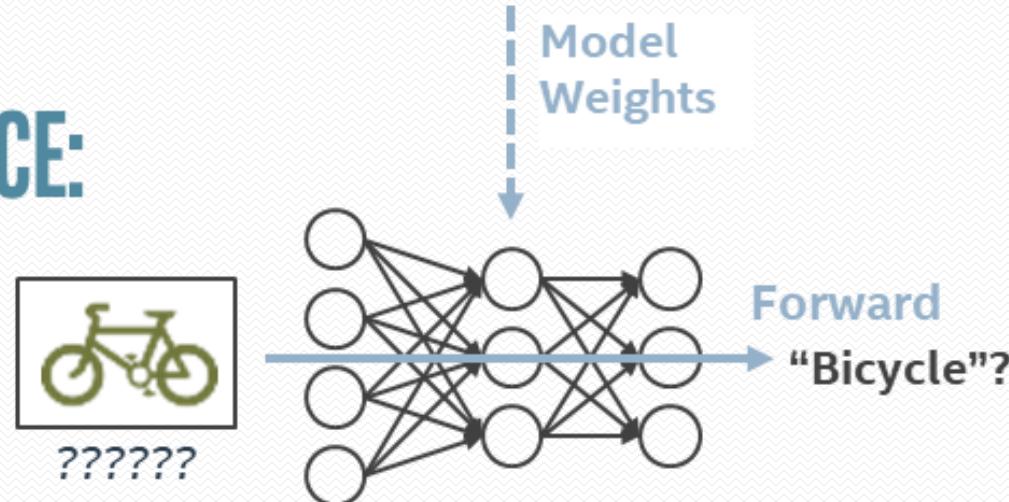
Technical

TRAINING & INFERENCE

TRAINING:



INFERENCE:

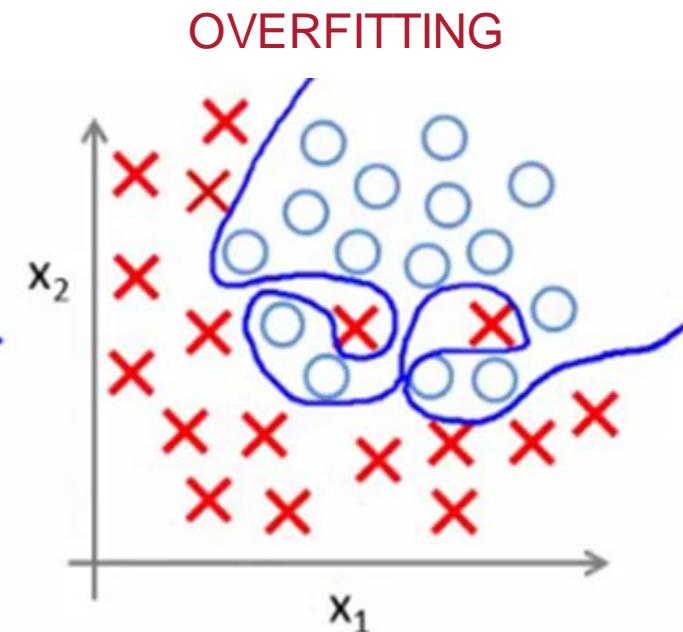
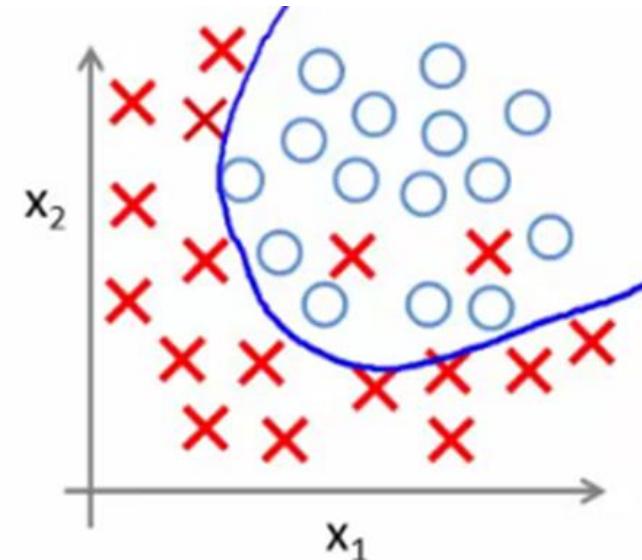
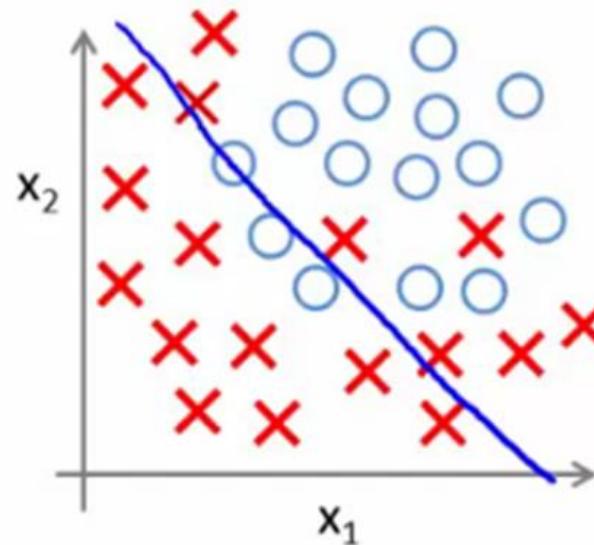


OVERFITTING / UNDERFITTING

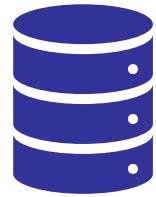
HIGH COMPLEXITY MODEL = HIGH POSSIBILITY OF OVERFITTING

If the model is too specific, generalization is reduced

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



DATASET



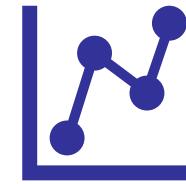
TRAIN

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

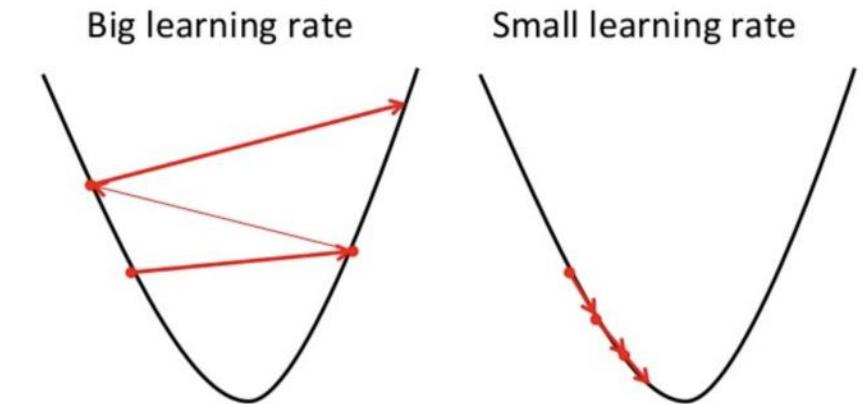
determines how fast or slow to move per iteration

INITIAL VALUE OF EACH PARAMETER

determines the initial value when running iteration

LOSS / COST FUNCTION

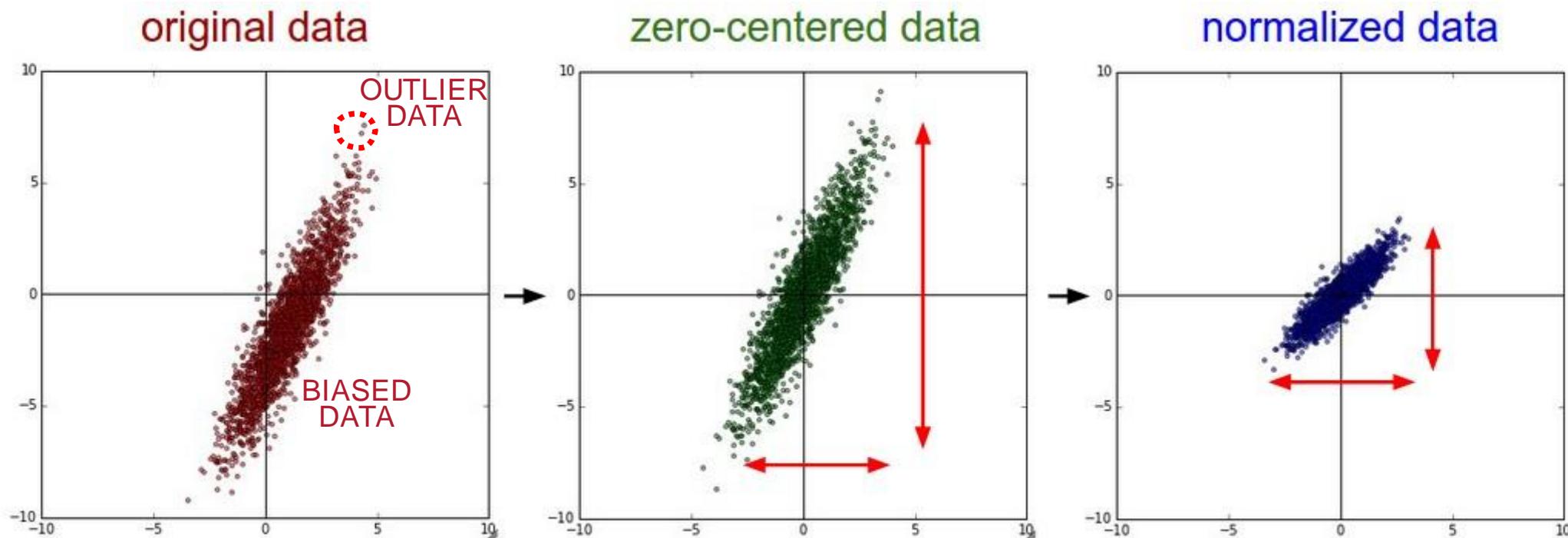
Decides how to calculate between actual result and calculated output



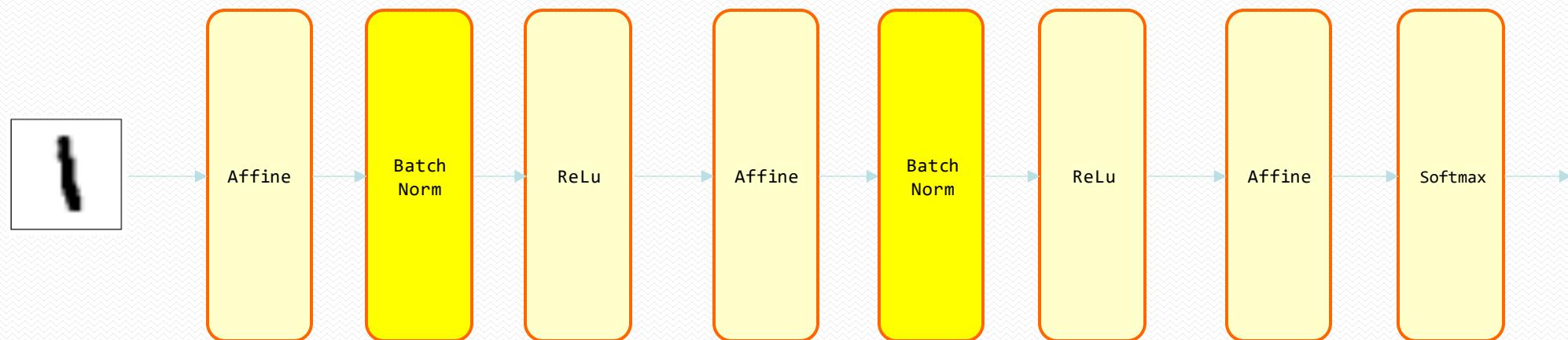
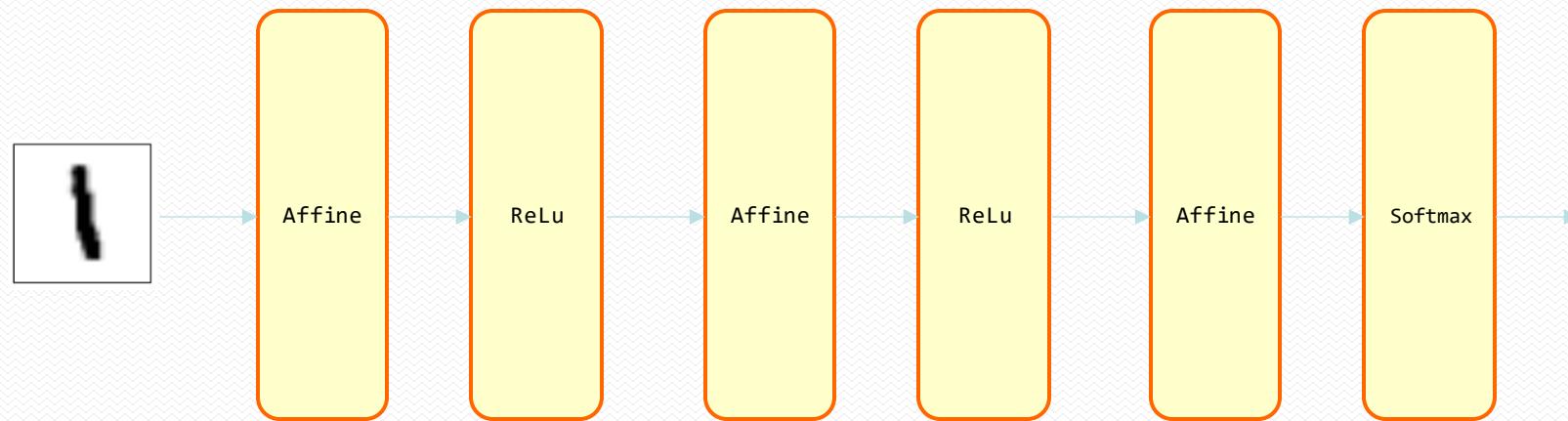
BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better



배치정규화 알고리즘

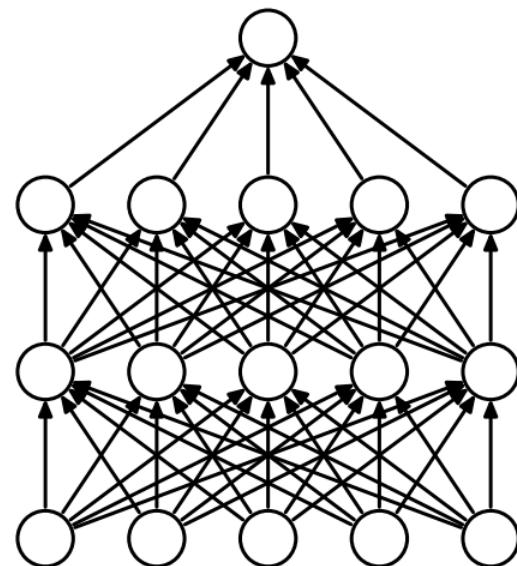


DROPOUT

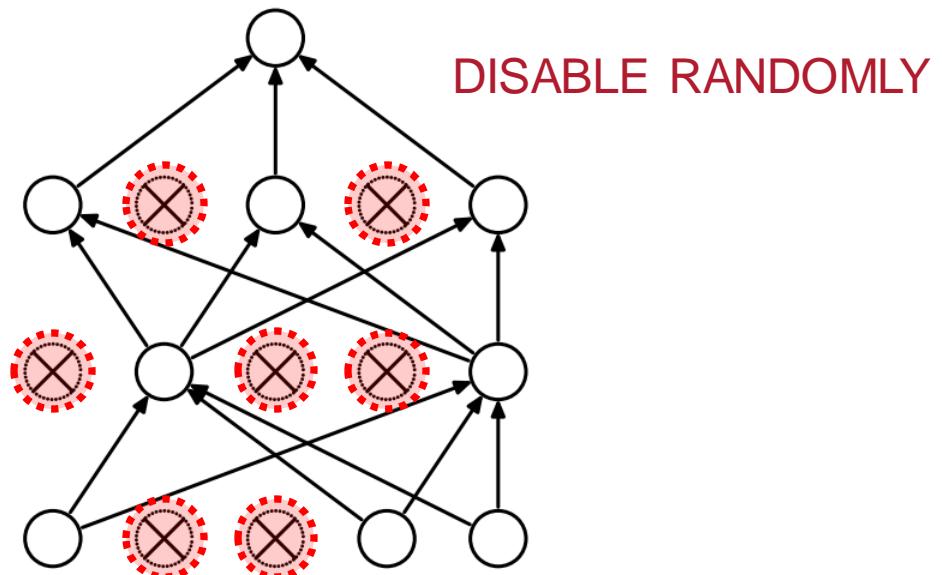
RANDOMLY SET SOME NEURONS TO ZERO IN THE FORWARD PASS

Voting effect → makes the result not biased

Anti co-adaptation → creates a robust network

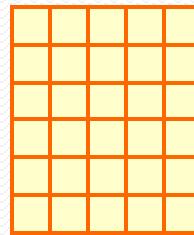
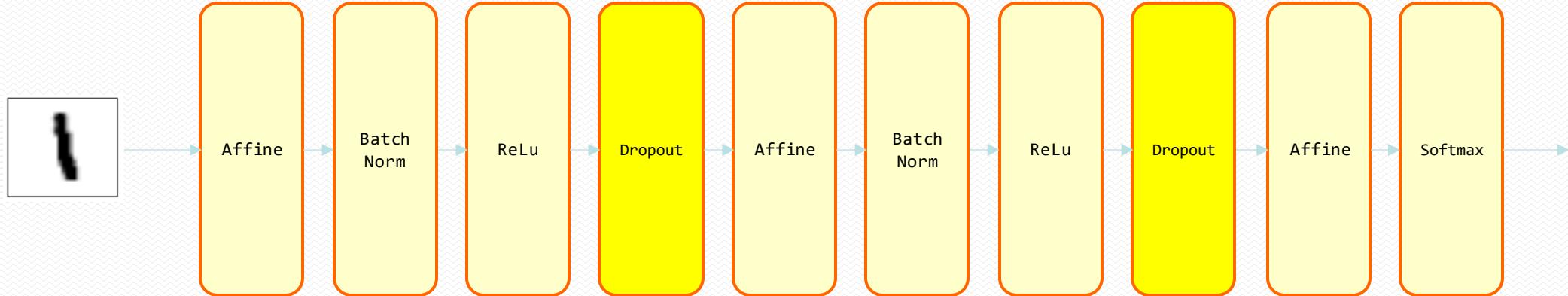


(a) Standard Neural Net

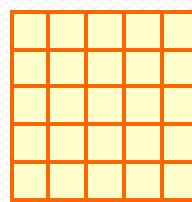


(b) After applying dropout.

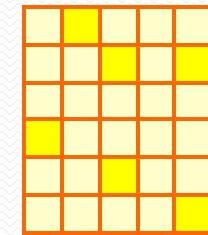
드롭아웃



(6, 5)

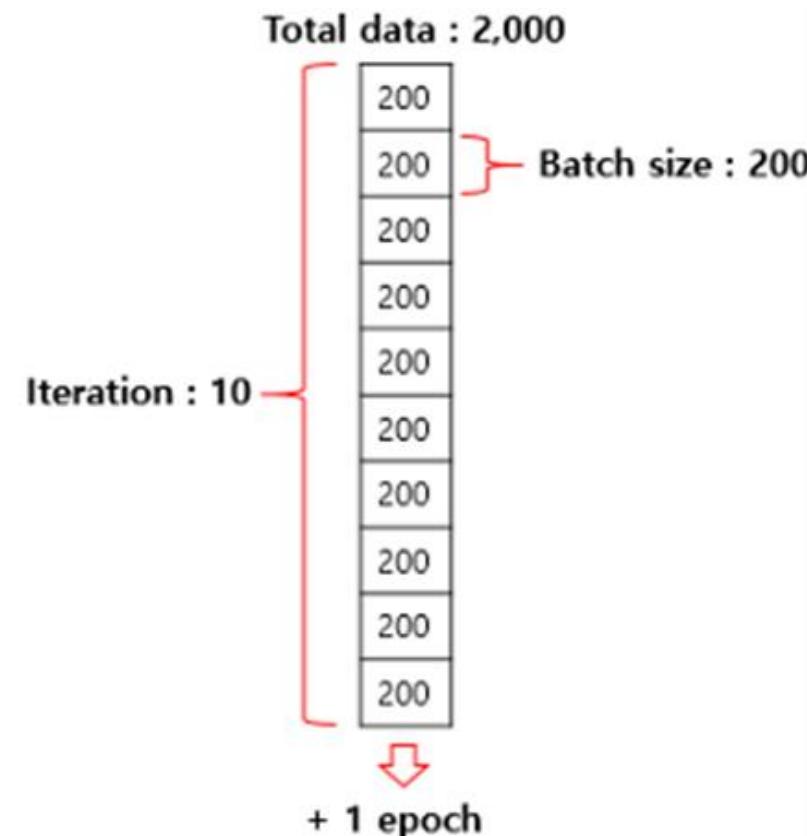


(5, 5)



(6, 5)

EPOCH & ITERATION & BATCH



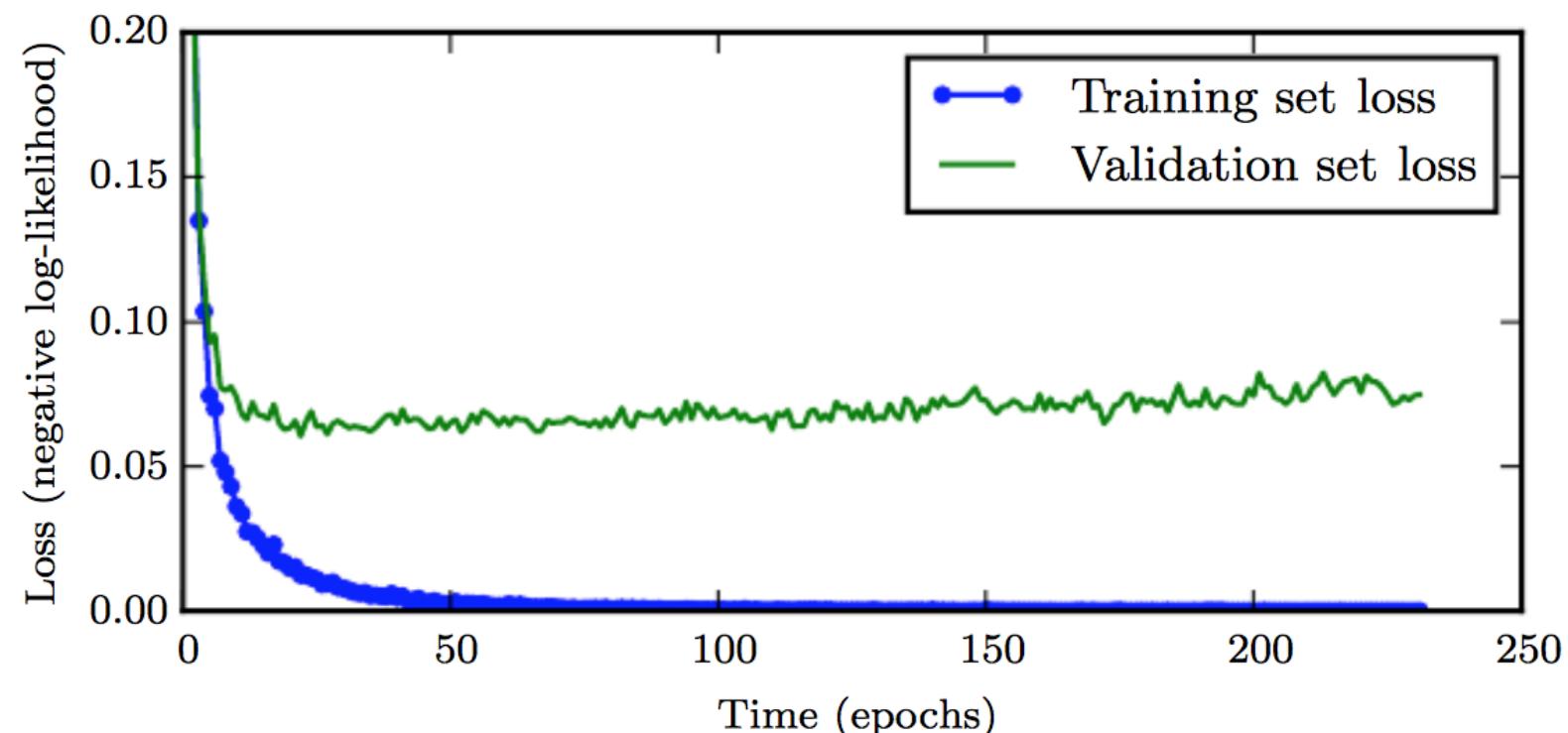
GENERALIZATION / RIZATION

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

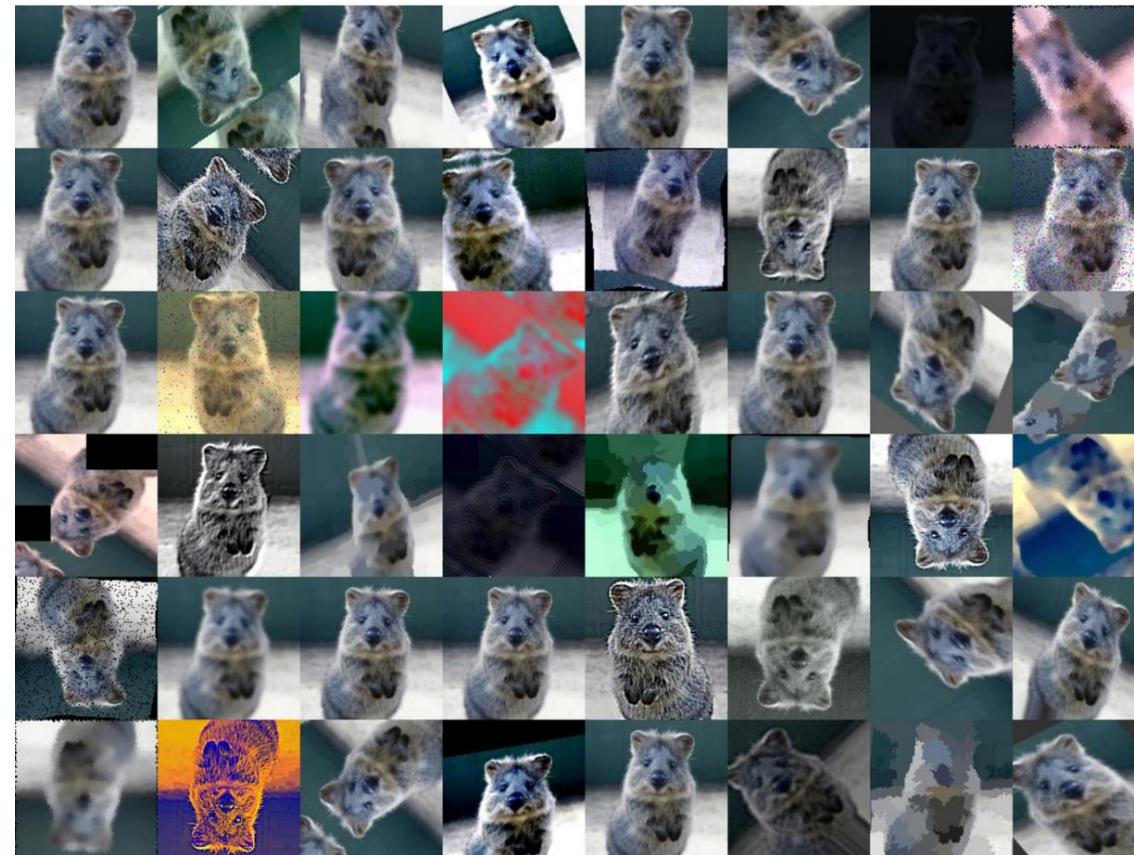
→ *UNDERFITTING AND OVERFITTING*

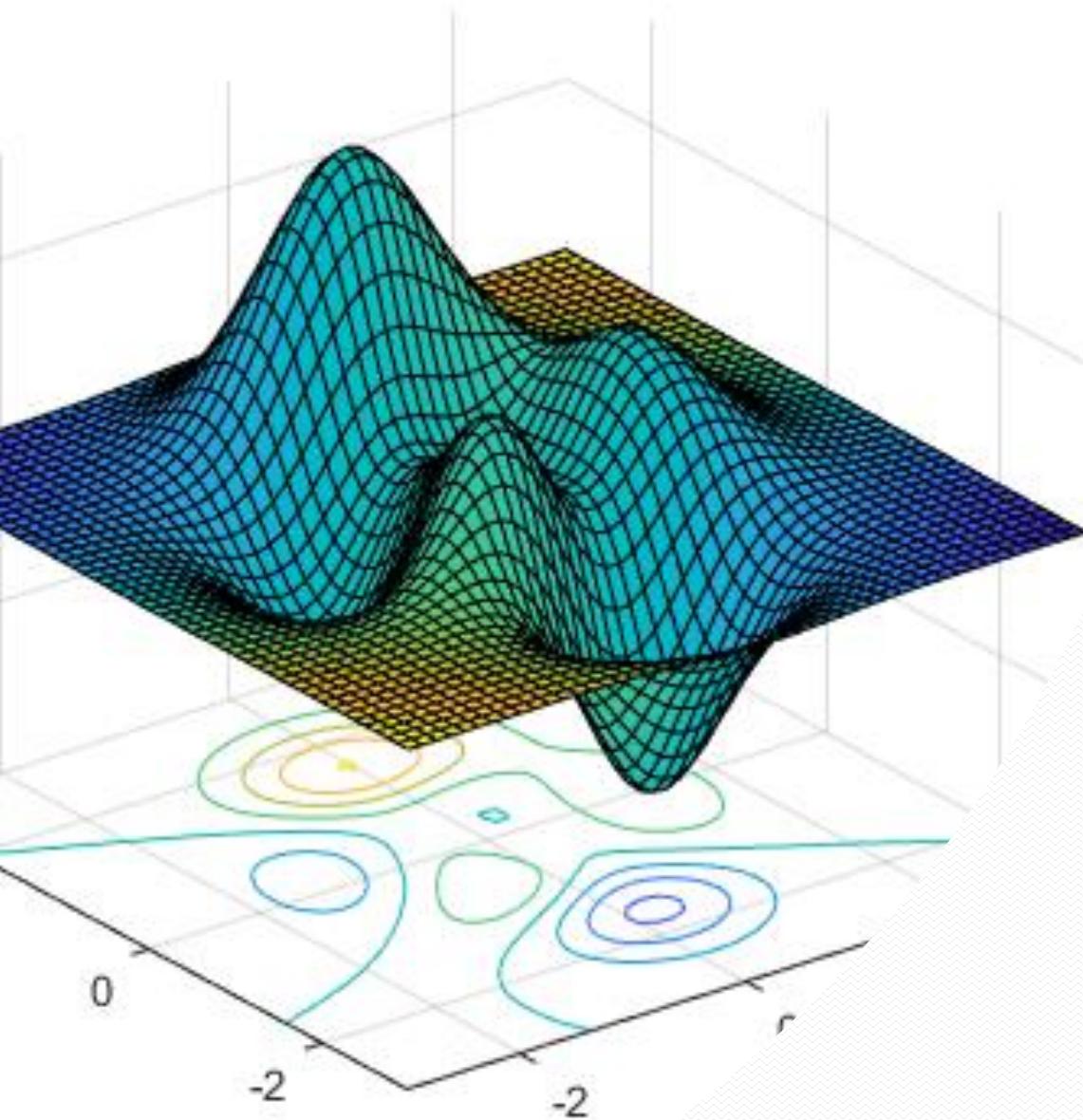


GENERALIZATION / REGULARIZATION

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





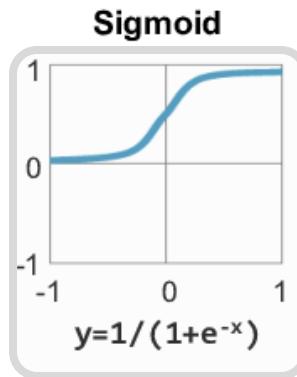
Deep-dive

Activation Function

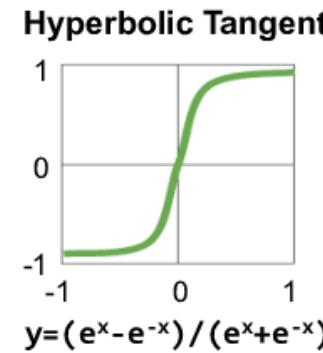
ACTIVATION FUNCTIONS

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

**Traditional
Non-Linear
Activation
Functions**

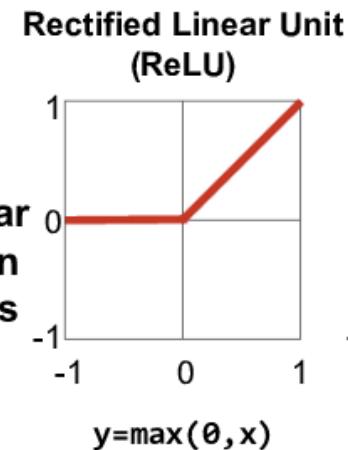


For binary classification

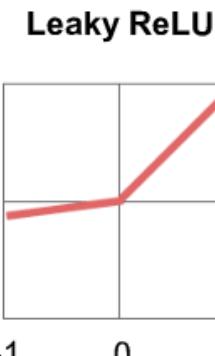


Widely used in hidden layers

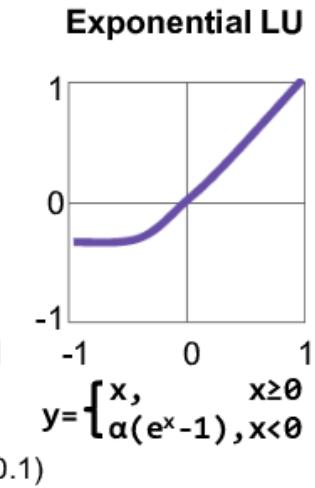
**Modern
Non-Linear
Activation
Functions**



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



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



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

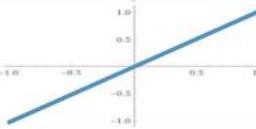
$$\alpha = \text{small const. (e.g. 0.1)}$$

ACTIVATION FUNCTIONS

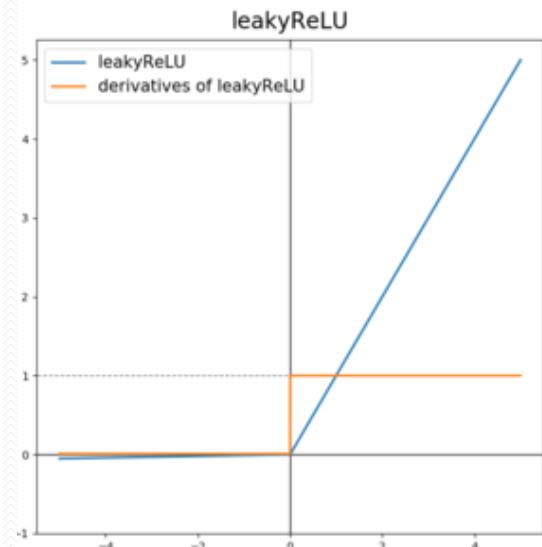
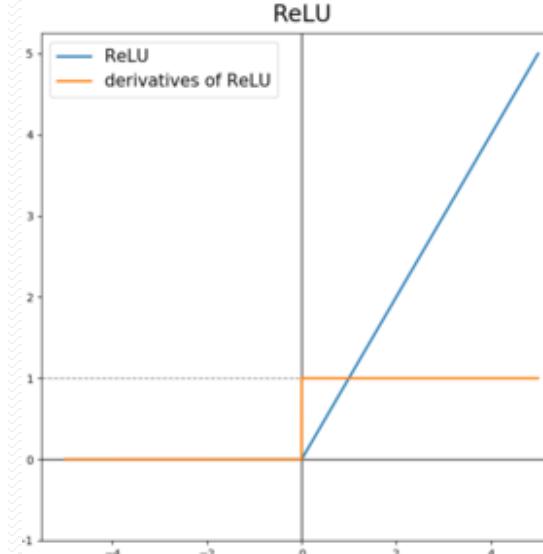
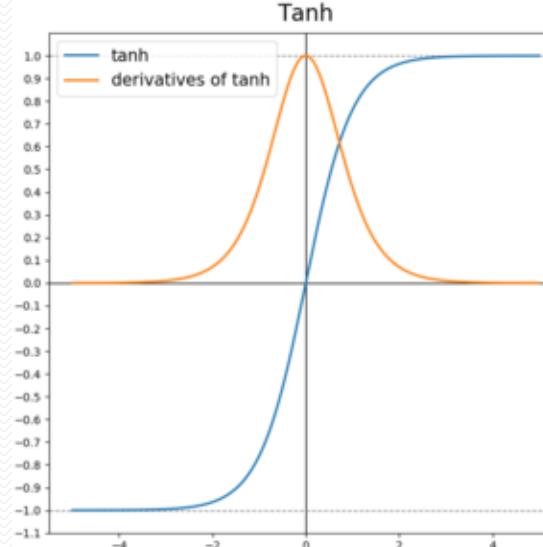
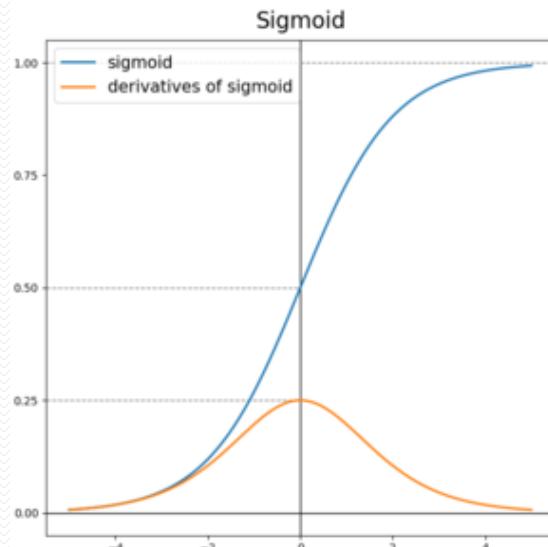
ACTIVATION FUNCTIONS

Why derivative/differentiation is used?

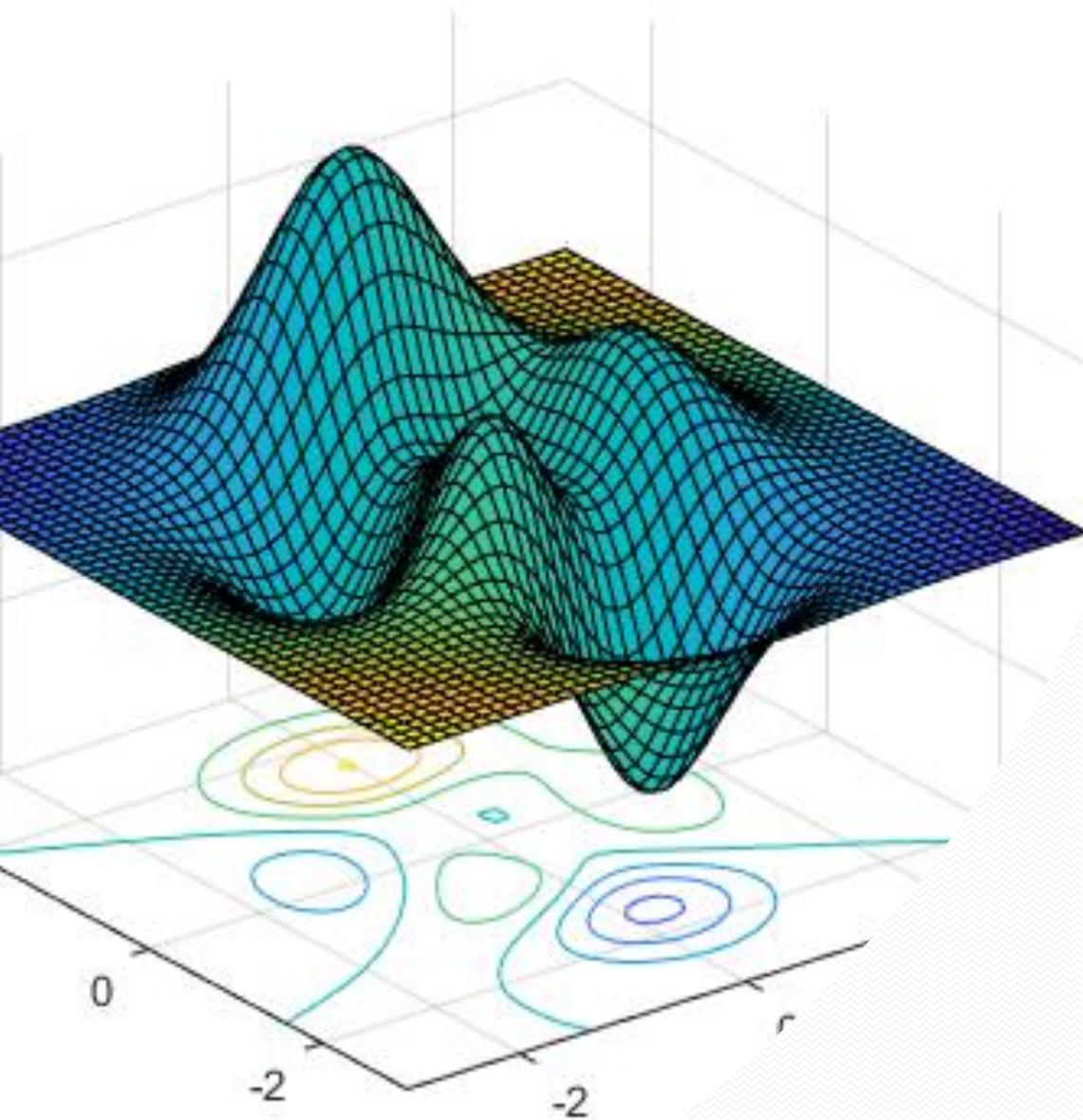
To know in which direction and how much to change or update the curve depending upon the slope.

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



Deep-dive

Optimizer

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.

Adadelta

- Extension of Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size.

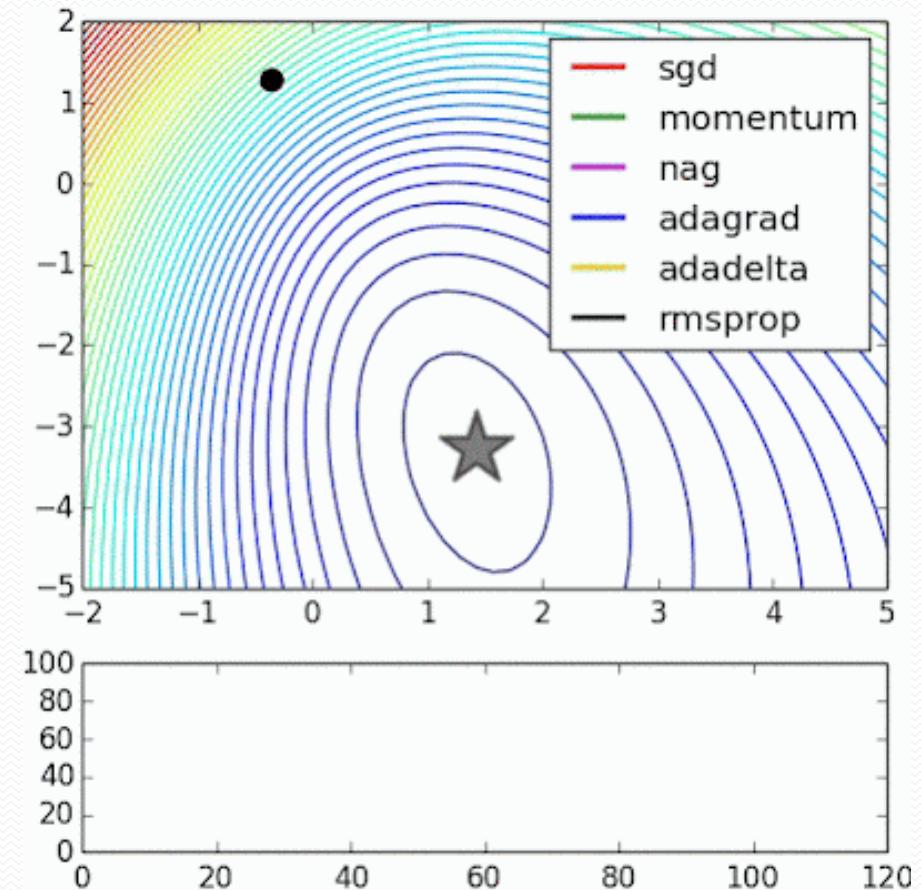
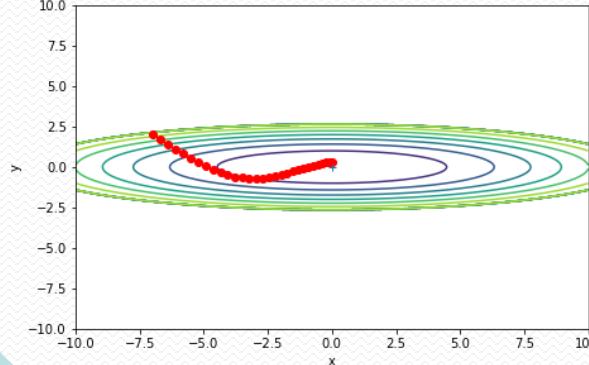
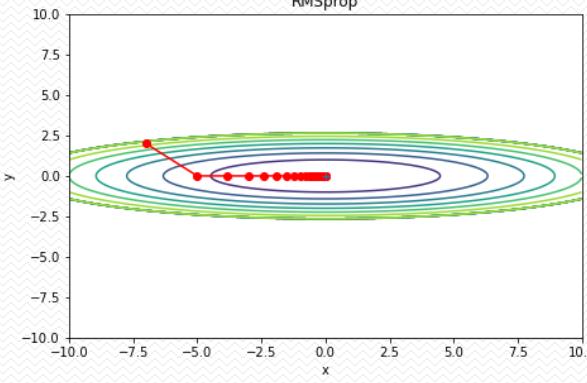
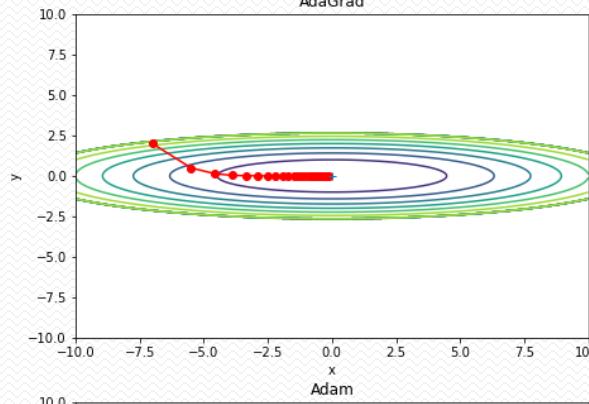
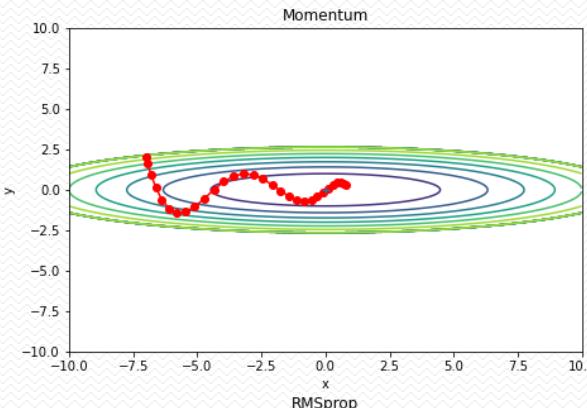
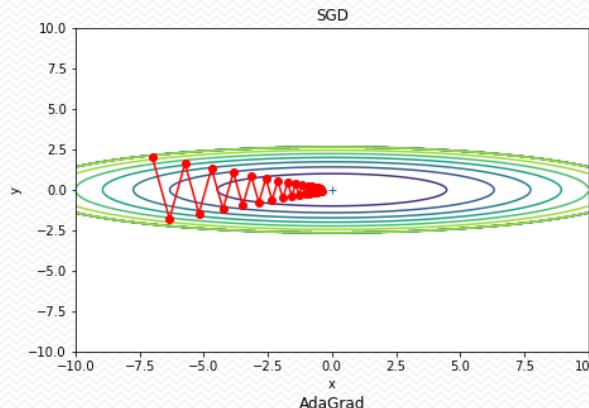
RMSProp

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

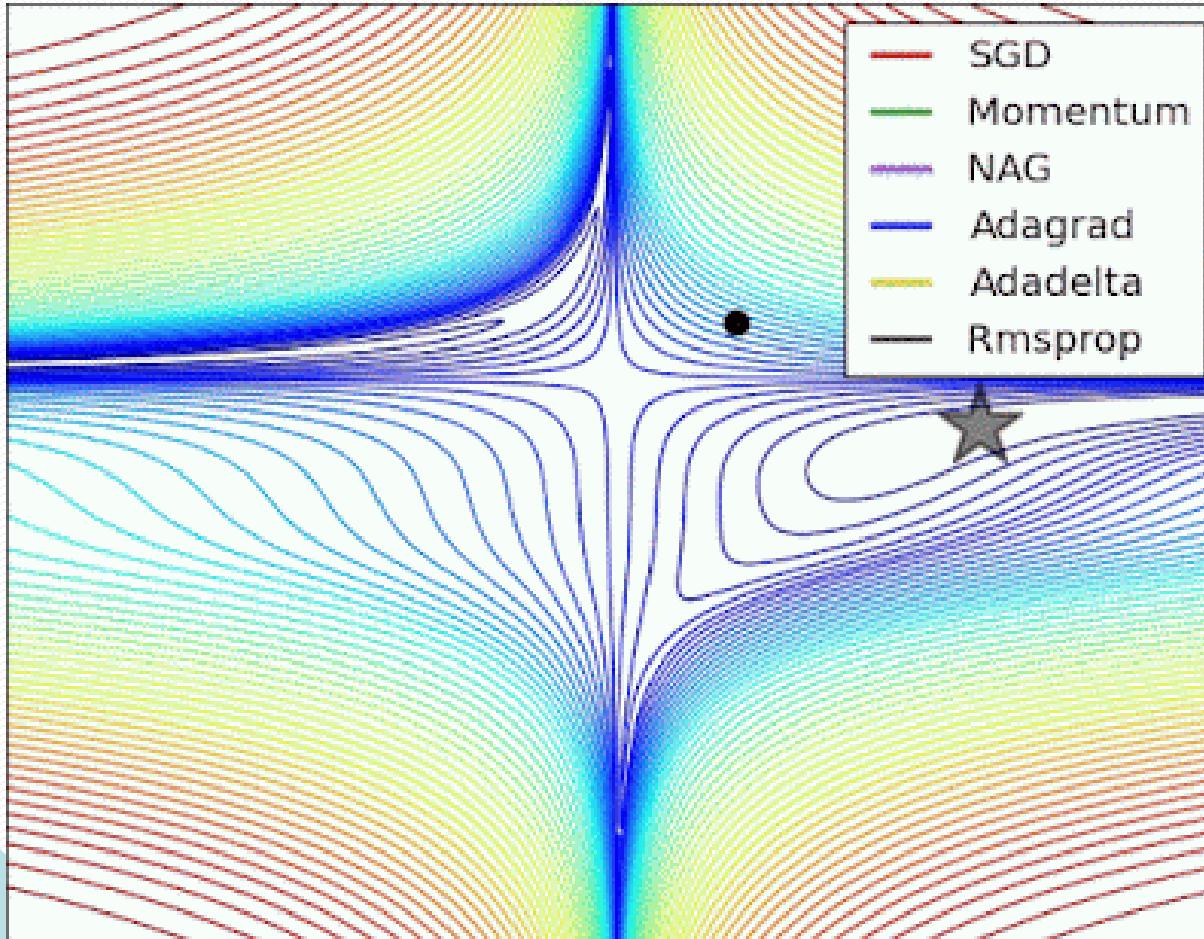
Adam (Adaptive Moment Estimation)

- 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

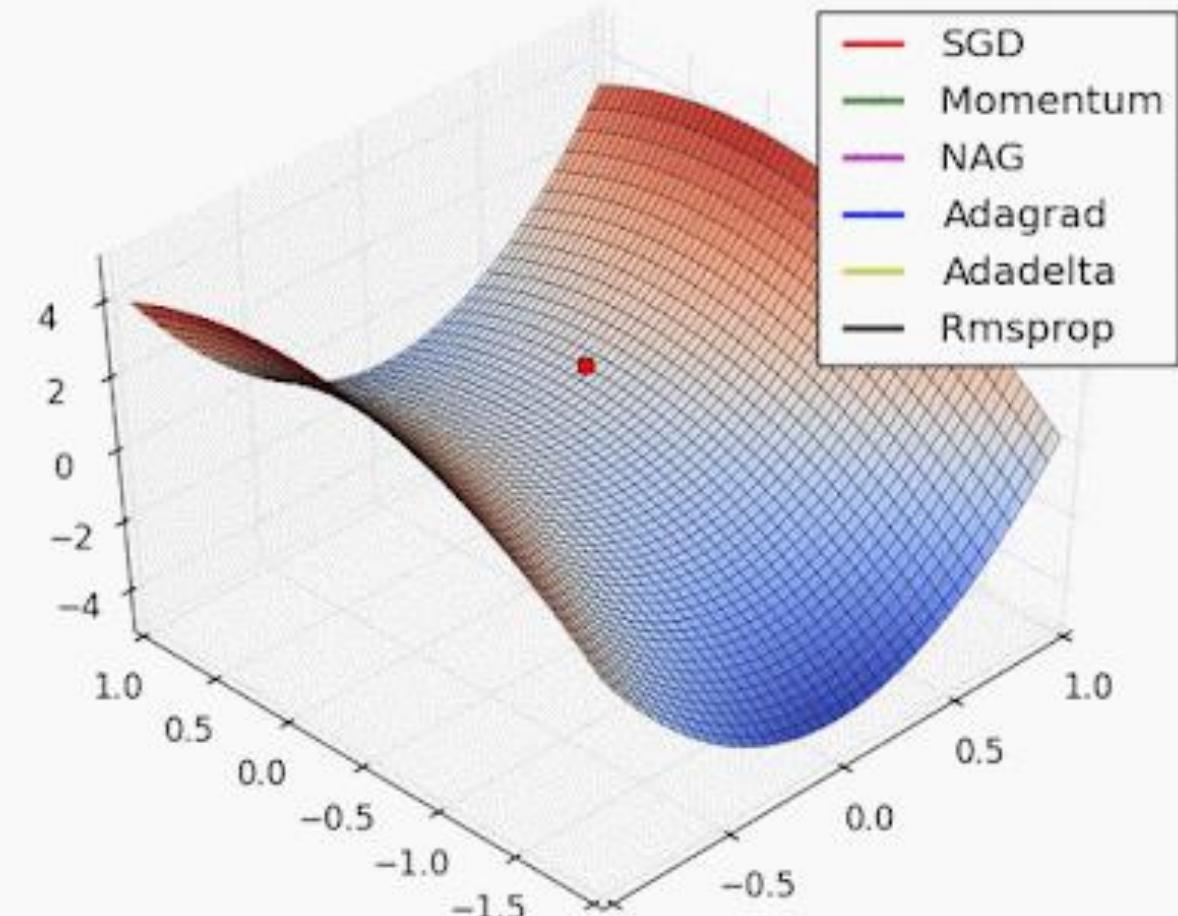
Behavior & Performance



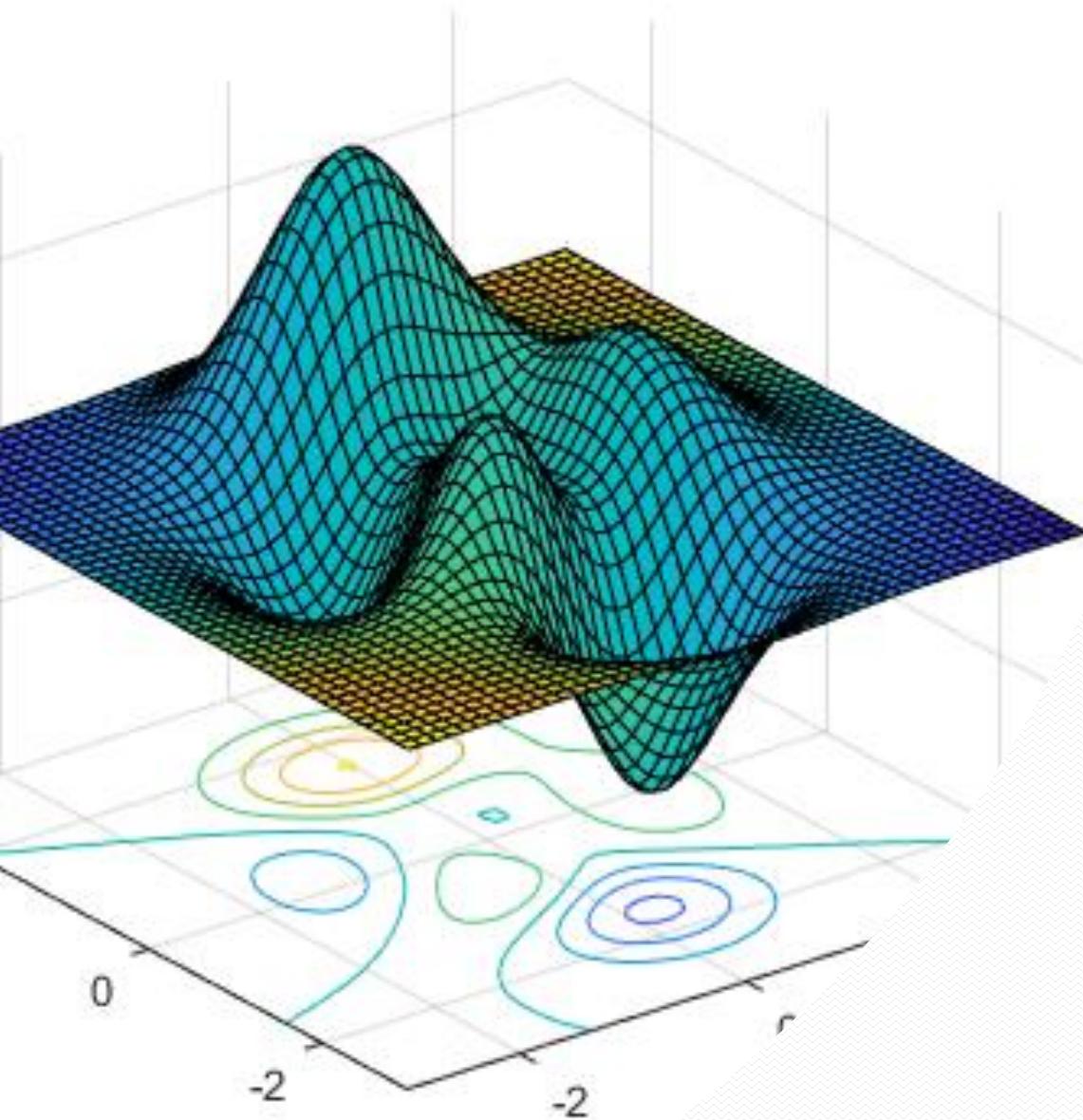
Behavior & Performance



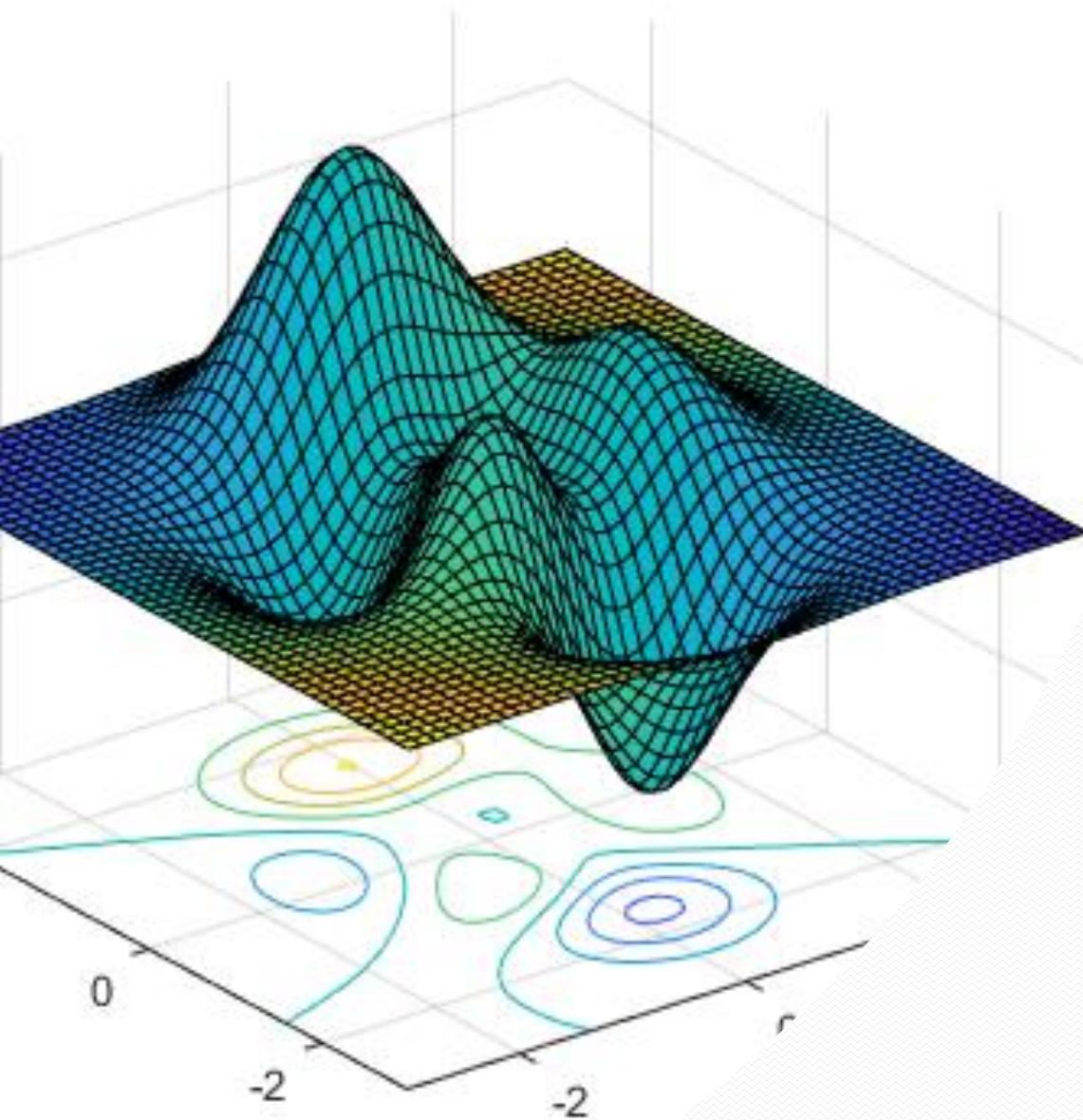
Saddle Point



Alec Radford has created some [great](#) [animations](#) comparing optimization algorithms



Hands-on



Back Propagation

MLP TRAINING & INFERENCE

FEEDFORWARD

Predict

BACKPROPAGATE

Weight Update

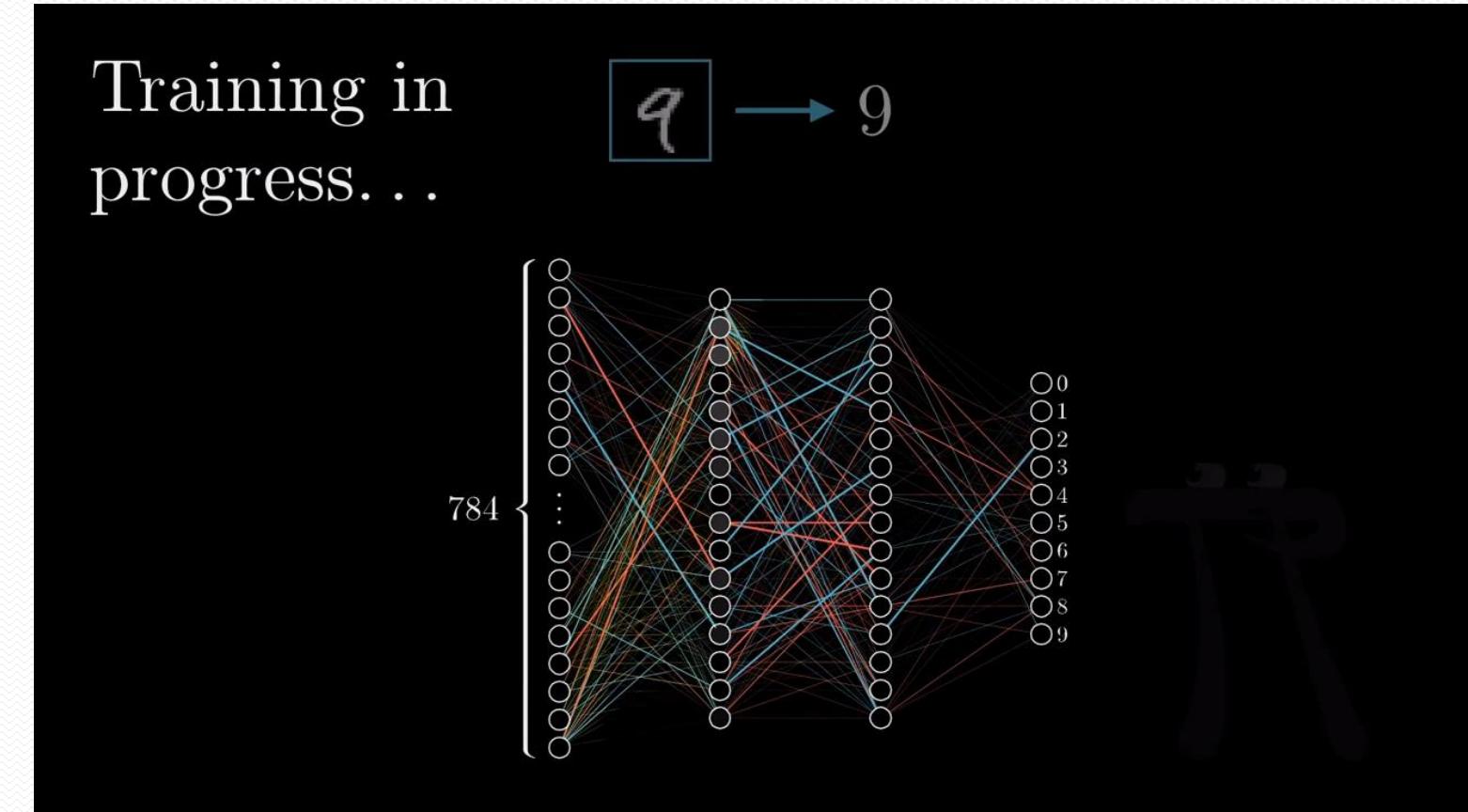
FEEDFORWARD

Predict

BACKPROPAGATE

Weight Update

...



CONTENTS

Session Break



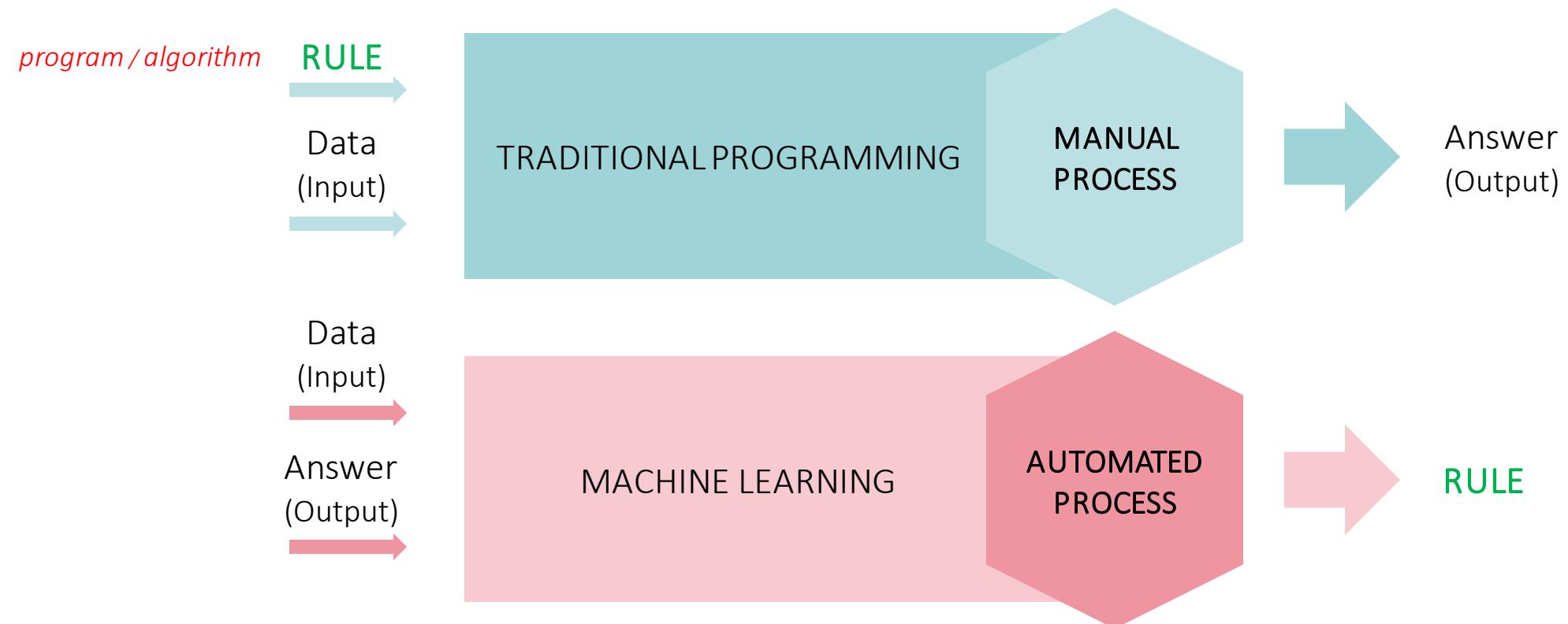
THANK YOU

Hands-on Colab

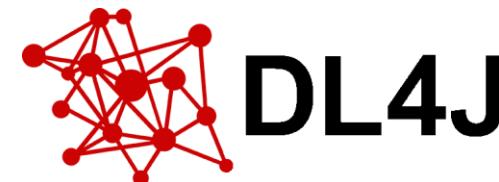
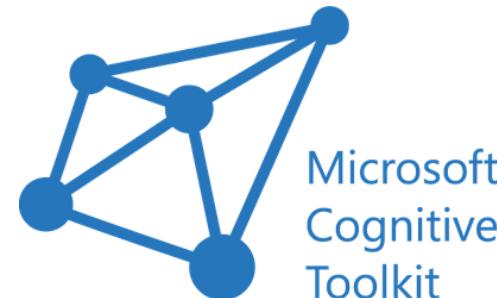
- [Hands-on] MNIST Dataset
 - https://colab.research.google.com/drive/122uCWPt1eR7yrasyG96ak0vXKKQGqa_Y?usp=sharing
- [Hands-on] Coco Dataset
 - https://colab.research.google.com/drive/1CrOUX7Ta-phwMI_Ngj1Ay9_DtyJUezLK?usp=sharing
- [Hands-on] ANN MNIST
 - https://colab.research.google.com/drive/1_ZhB7hwtYCEtfHTwkUthEsLHx3rlehjz?usp=sharing
- [Hands-on] Benchmark app
 - https://colab.research.google.com/drive/1mF99L-U5NJOKYjf_VY2ZkA6vhcpX7CSh?usp=sharing

TRADITIONAL PROGRAMMING / ML

Traditional Programming refers to any manually created program
uses input data and runs on a computer to produce the output



DEEP LEARNING FRAMEWORKS



SUMMARY

AI SIMULATING LIVING BEHAVIOR

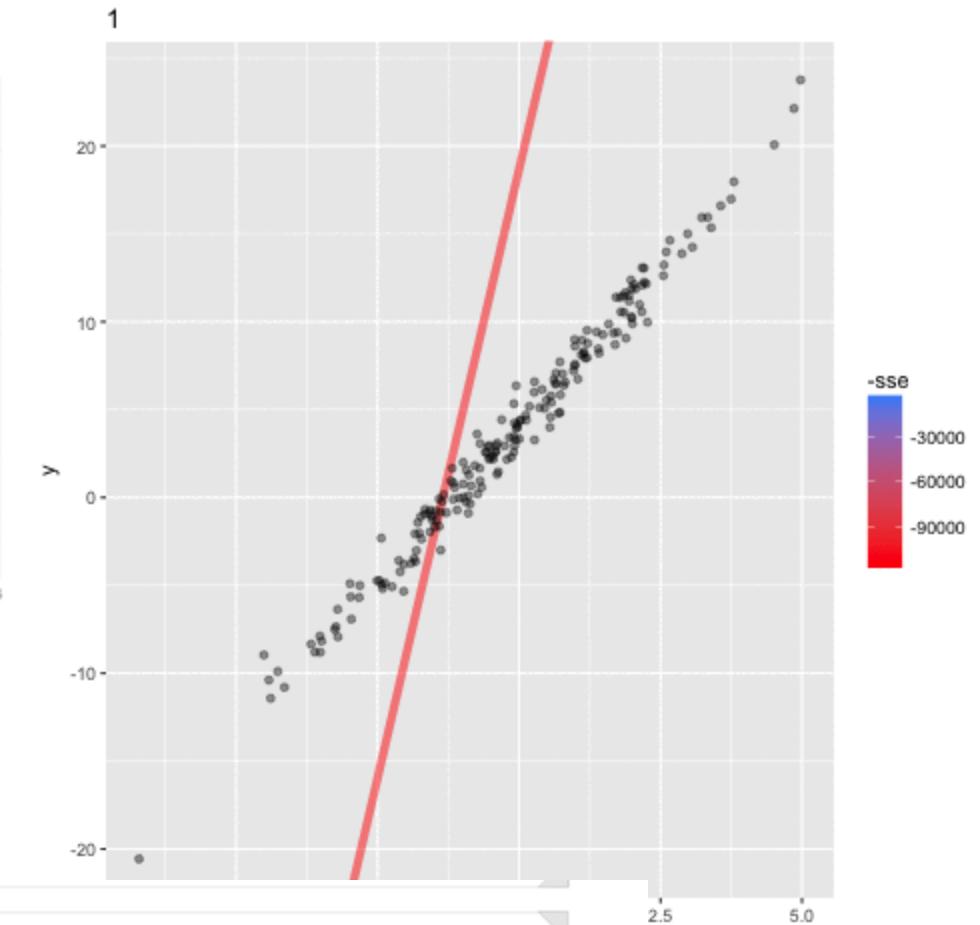
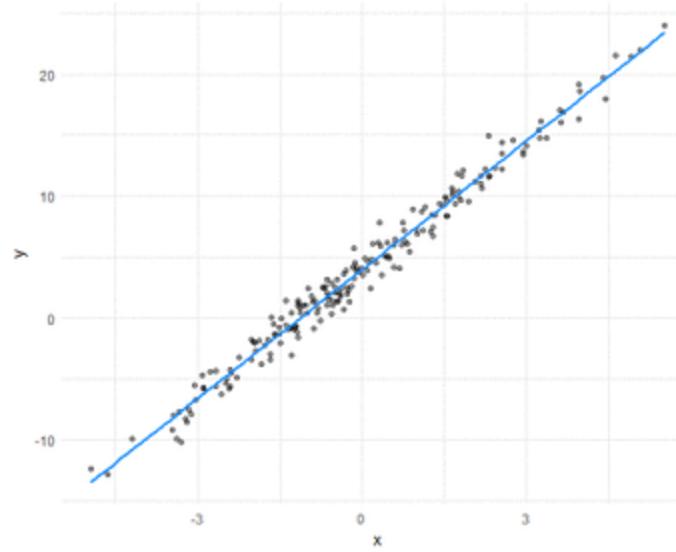
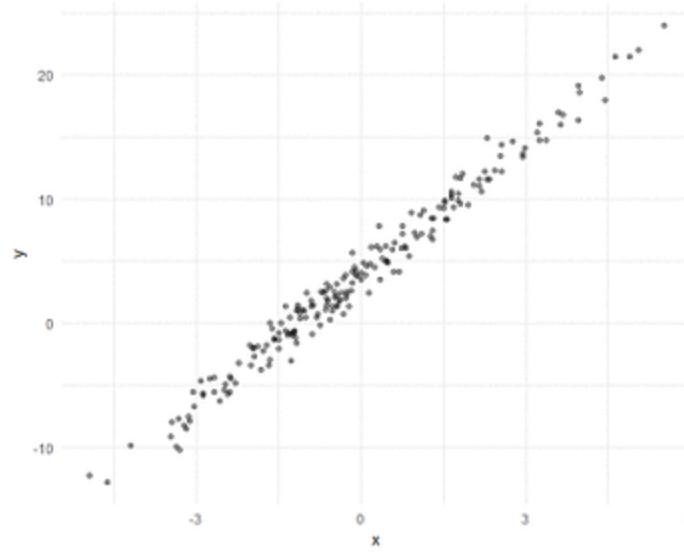
ML TRAINING MODEL FROM DATASETS

DL INSPIRED BY THE NEURON NETWORKS IN BRAIN

FEATURE EXTRACTION , DL

No need to manually extract features from the image
Layered feature extraction

Lowest cost



Hypothesis :

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

Parameters:

$$w, b$$

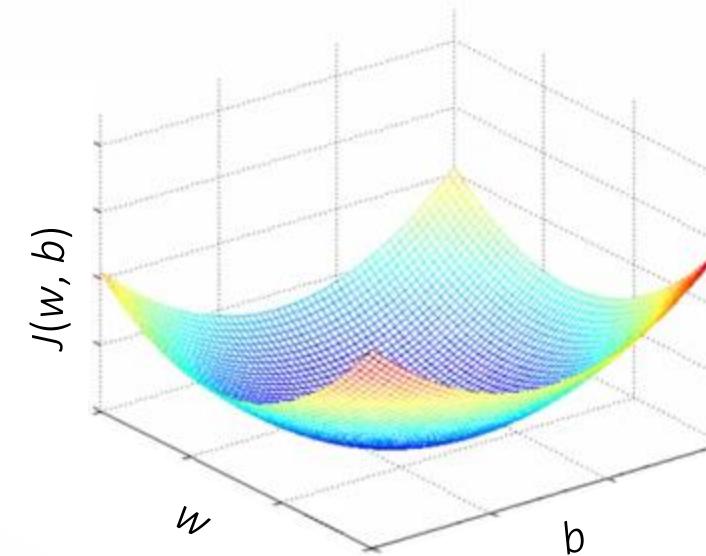
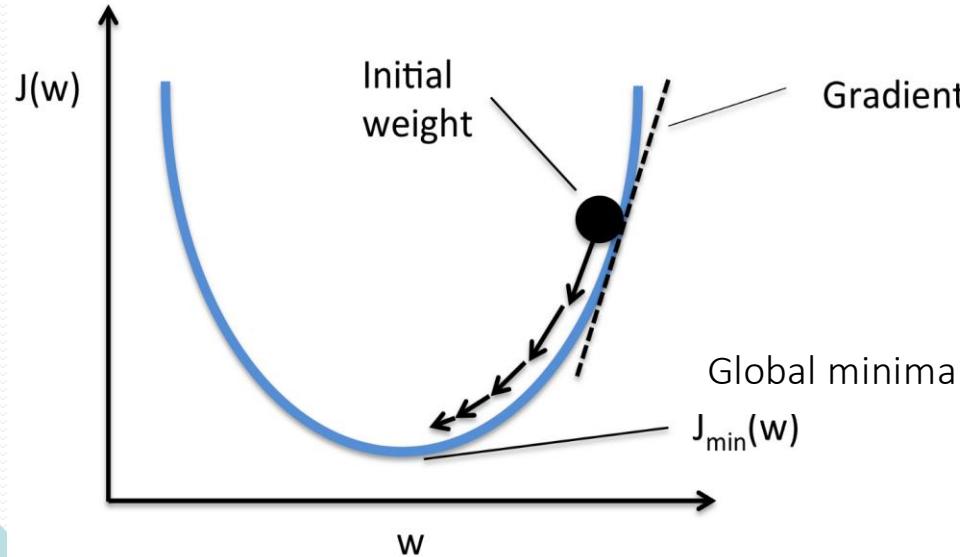
Cost(Loss) Function

$$Cost(w, b) = \frac{1}{2}(\hat{y} - y)^2$$

$$Cost = \frac{1}{n} \sum_{1}^n (H(x_n) - y_n)^2$$

Gradient Descent

- Example: mean squared error as loss function
- Cost function is a convex function
- **Goal:** Find values of w and b such that J is minimum i.e. find the global minima of the cost function
- Move in the direction of negative gradient



$$w' = w - \alpha \frac{\partial J}{\partial w}$$

learning rate

$$b' = b - \alpha \frac{\partial J}{\partial b}$$