



PyTorch 배우는 딥러닝 입문 (1)

Introduction to Deep Learning

이영완

한국전자통신연구원



이영완

ETRI 초지능창의연구소 시각지능연구실 (2017 ~)
Ph.D candidate in Graduate school of AI @KAIST

- 주요 연구 분야 ([Google scholar](#))
 - Generative model
 - KOALA (NeurIPS`24)
 - Self-supervised learning
 - EVEREST (ICML`24)
 - RC-MAE (ICLR`23)
 - Neural architecture design:
 - VoVNet (CVPRW`19, CVPR`20)
 - MPViT (CVPR`22)
 - Object detection and segmentation
 - CenterMask (CVPR`20)
 - Video classification
 - VoV3D (IEEE Access)
- 수상 이력
 - 장관표창, 과학기술정보통신부 2022
 - 우수논문상, ETRI 2022, 2024
 - 우수강사상, ETRI 2022
 - #2 rank ImageNet challenge (detection task, 2017)



Youngwan Lee

Visual Intelligence lab
@ETRI

Senior researcher at ETRI &
Ph.D student at KAIST

📍 Daejeon, South Korea

✉ Email

LinkedIn

DBLP

Github

Google Scholar

ORCID

I'm Youngwan, a senior researcher at [ETRI](#) and Ph.D student in Graduate school of [AI](#) at [KAIST](#), where I'm advised by Prof. [Sung Ju Hwang](#) in the [Machine Learning and Artificial Intelligence \(MLAI\)](#) lab. My research interest is how computers understand the world, including efficient 2D/3D neural network design, object detection, instance segmentation, semantic segmentation, and video classification . Recently, I have focused on Vision Transformer architecture and self-supervised learning.

News

- 2024.09: 🎉 KOALA has been accepted at NeurIPS 2024.
- 2024.08: 🎉 DiT-Pruner has been accepted at ECCVW 2024.
- 2024.05: 🎉 EVEREST has been accepted at ICML 2024.
- 2023.12: 🎉 KOALA, a fast Text-to-Image synthesis model, has been released.
- 2023.02: 🎉 Two papers have been accepted at ICLR 2023.
- 2022.02: 🎉 One paper has been accepted at CVPR 2022.

Publications

[Google Scholar full list](#)

• KOALA: Empirical Lessons Toward Memory-Efficient and Fast Diffusion Models for Text-to-Image Synthesis

Youngwan Lee, Kwanyong Park, Yoorhim Cho, Yong-Ju Lee, Sung Ju Hwang

Advances in Neural Information Processing Systems ([NeurIPS](#)) 2024

CVPR 2024 Workshop on [Generative Models for Computer Vision](#)

[\[Project page\]](#)[\[paper\]](#)[\[code\]](#)

• DiT-Pruner: Pruning Diffusion Transformer Models for Text-to-Image Synthesis Using Human Preference Scores

Youngwan Lee, Yong-Ju Lee, Sung Ju Hwang

European Conference on Computer Vision ([ECCV](#)) 2024 Workshop on [Green Foundation Models](#)

[\[paper\]](#)



Youngwan Lee

youngwanLEE

Senior researcher at ETRI & Ph.D
student in Graduate school of AI at
KAIST.

Edit profile

281 followers · 3 following

ETRI & MLAI@KAIST

South Korea

21:24 (UTC +09:00)

yw.lee@etri.re.kr

<https://youngwanlee.github.io/>

in/youngwanlee

<https://github.com/youngwanLEE>

<https://scholar.google.com/citations?user=EqemKYsAAAAJ&hl=ko>

[sdxl-koala](#) Public PyTorch

Compressing SDXL via knowledge-distillation

⭐ 125 ⚡ 3

[rc-mae](#) Public PyTorch

[ICLR 2023] RC-MAE

Python ⭐ 49 ⚡ 2

[MPViT](#) Public PyTorch

[CVPR 2022] MPViT:Multi-Path Vision Transformer for Dense Prediction

Python ⭐ 354 ⚡ 39

[UAD](#) Public PyTorch

[ECCVW 2022] UAD: Localization Uncertainty Estimation for Anchor-Free Object Detection

Python ⭐ 13

[centermask2](#) Public PyTorch

[CVPR 2020] CenterMask : Real-time Anchor-Free Instance Segmentation

Python ⭐ 771 ⚡ 159

[vovnet-detectron2](#) Public PyTorch

[CVPR 2020] VoVNet backbone networks for detectron2

Python ⭐ 372 ⚡ 69



Youngwan Lee

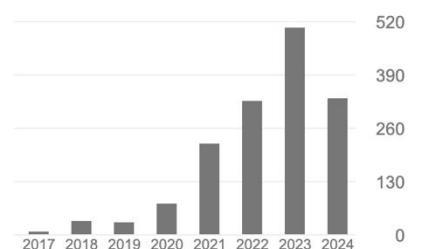
[Electronics and Telecommunications Research Institute; KAIST](#)
etri.re.kr의 이메일 확인됨 - [홈페이지](#)

Computer Vision Machine Learning

팔로우

인용

	전체	2019년 이후
서지정보	1540	1497
h-index	11	10
i10-index	13	10



공동 저자

수정

- | 제목 | 인용 | 연도 |
|--|-----|------|
| CenterMask: Real-time anchor-free instance segmentation
Y Lee, J Park
Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern ... | 631 | 2020 |
| An Energy and GPU-Computation Efficient Backbone Network for Real-Time Object Detection
Y Lee, J Hwang, S Lee, Y Bae, J Park
Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern ... | 391 | 2019 |
| MPViT: Multi-Path Vision Transformer for Dense Prediction
Y Lee, J Kim, J Willette, SJ Hwang
Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern ... | 243 | 2022 |

강의 목표

- '알아두면 쓸모 있는 신박한 딥러닝 입문' 탑재해드립니다.
- 딥러닝 오픈소스 대세인 Pytorch 첫 출발을 도와드립니다.
- Pytorch를 이용한 연구 개발 경험과 꿀팁을 드립니다.
- 자유롭게 질의응답할 수 있는 편안한 토론의장을 선사합니다.

강의 추천 대상자

- 딥러닝에 대해 기초부터 차근차근 배우고 싶은 분
여기저기서 딥러닝 딥러닝하는데 도대체 딥러닝이 무엇인지 궁금하신 분
- 딥러닝을 시작하고 싶은데 어떻게 시작해야 될지 모르는 분
더 늦기 전에 딥러닝을 시작해야 될텐데... 라고 생각하시는 분
- 코드로 딥러닝 개념을 이해하고 싶은 분
글이나 말보다는 코드로 세상을 이해하시는 분
- Pytorch를 업무에 활용하고 싶은 분
요즘 딥러닝은 pytorch가 대세라던데... pytorch에 밭을 담그고자 하시는 분

강의 내용

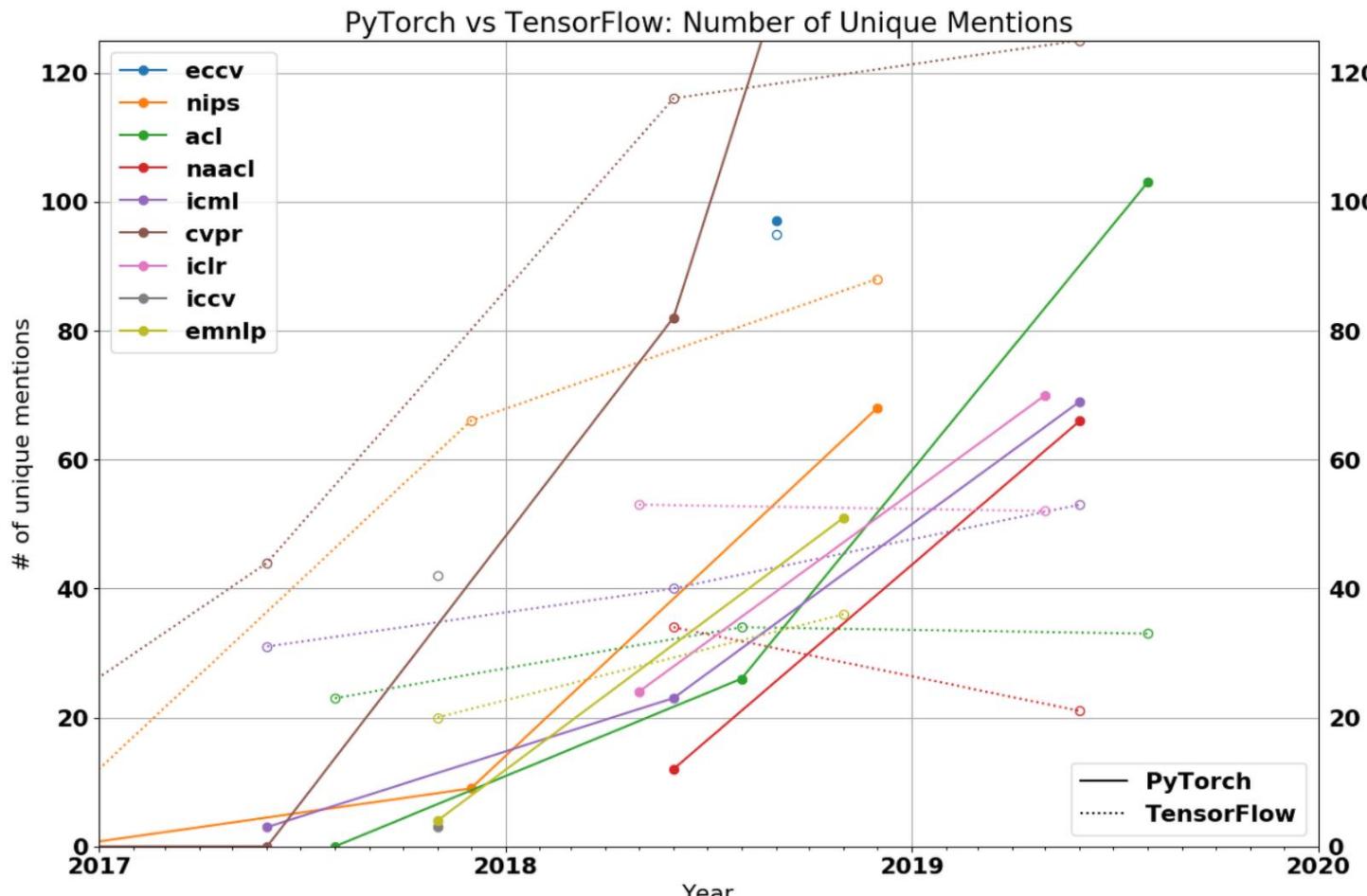
[1day]

- 이론 : 딥러닝 기초
 - Overview of Machine Learning (ML)
 - Introduction to Neural Network
- 실습 : Pytorch Tutorial & AutoGrad

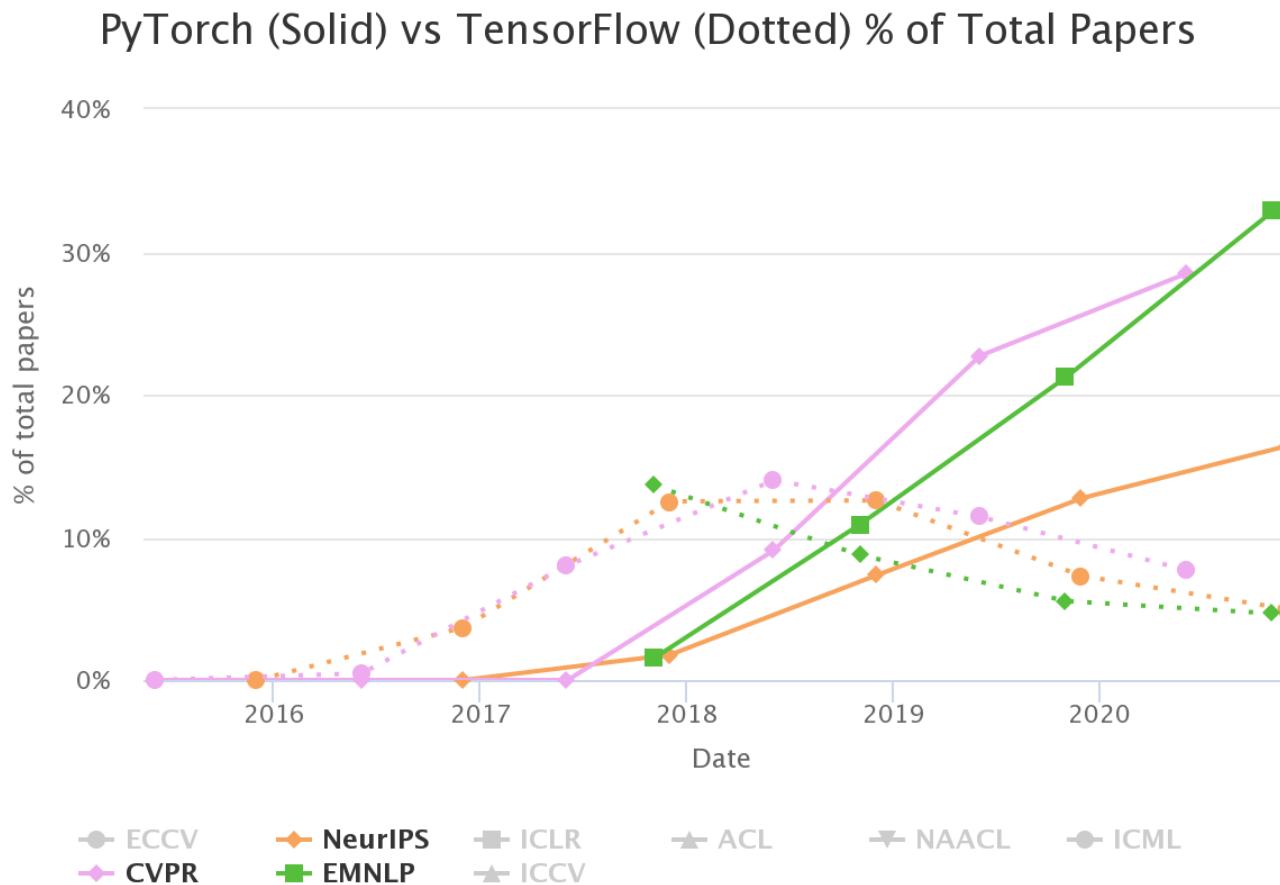
[2day]

- 이론 : Convolutional Neural Network (CNN)
- 실습 : CNN 네트워크 구현 및 Image classification 실습

Why Pytorch?



Why Pytorch?



Artificial Intelligence in Context of Human History

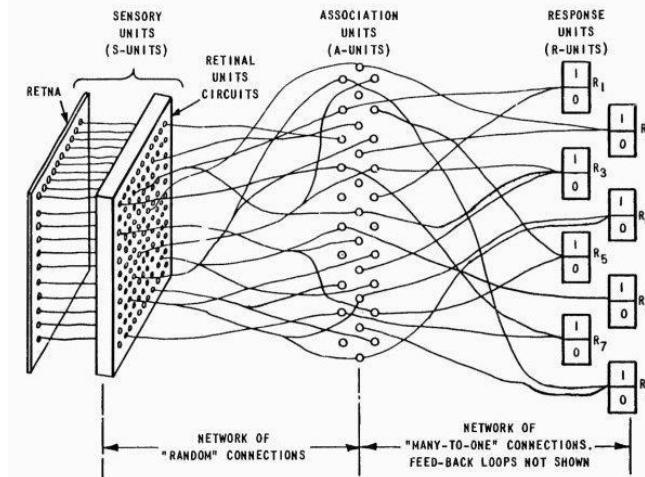
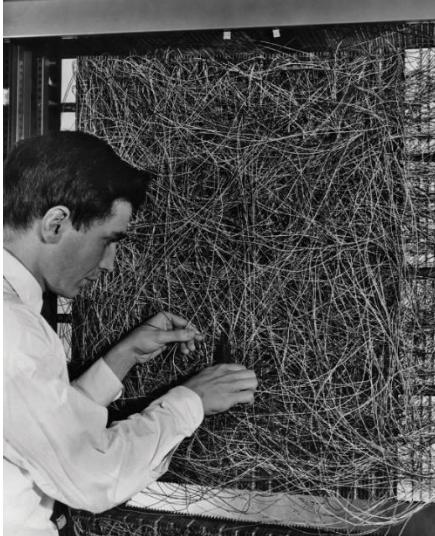


Figure 1 ORGANIZATION OF THE MARK I PERCEPTRON

Dreams, mathematical foundations, and engineering in reality.

Frank Rosenblatt, Perceptron (1957, 1962): Early description and engineering of singlelayer and multilayer artificial neural network

Artificial Intelligence in Context of Human History



Kasparov vs Deep Blue, 1997

Artificial Intelligence in Context of Human History



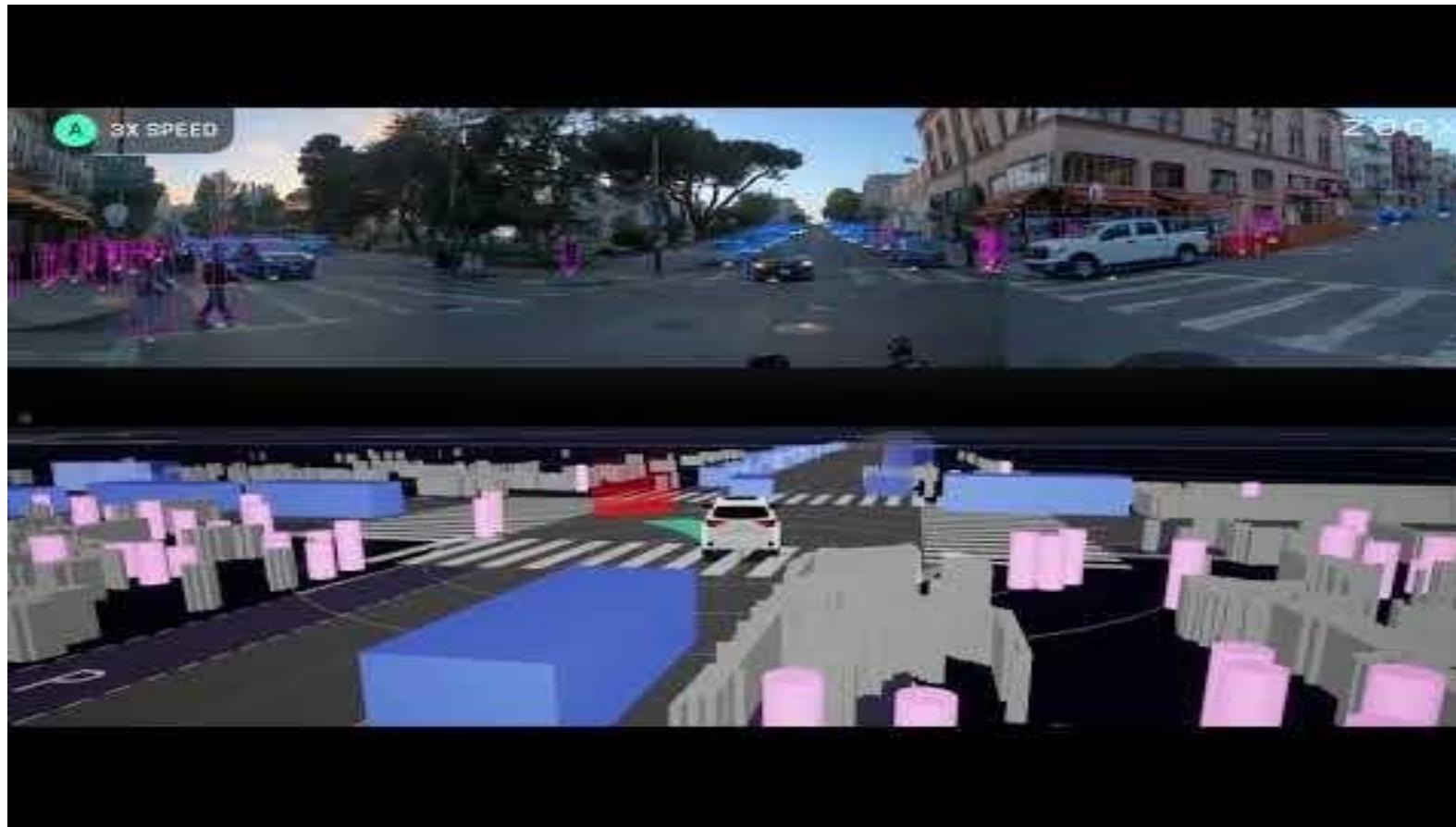
Lee Sedol vs AlphaGo, 2016

Artificial Intelligence in Context of Human History



Self-driving car

Artificial Intelligence in Context of Human History



Visual perception for Self-driving car

Artificial Intelligence in Context of Human History



Robotics

Artificial Intelligence in Context of Human History

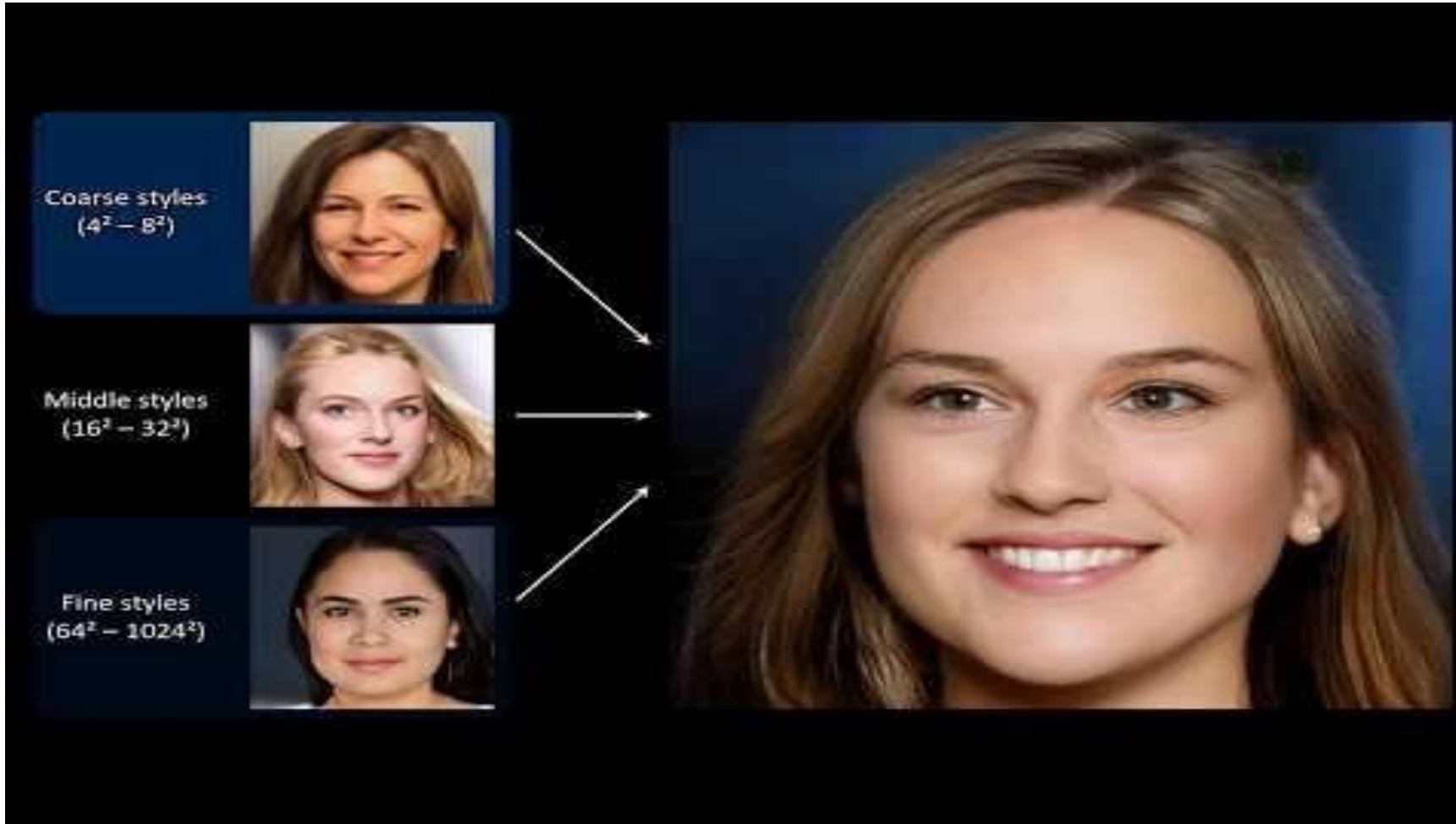


Image generation

Artificial Intelligence in Context of Human History



Image generation

Artificial Intelligence in Context of Human History



Audio style transfer

Artificial Intelligence in Context of Human History



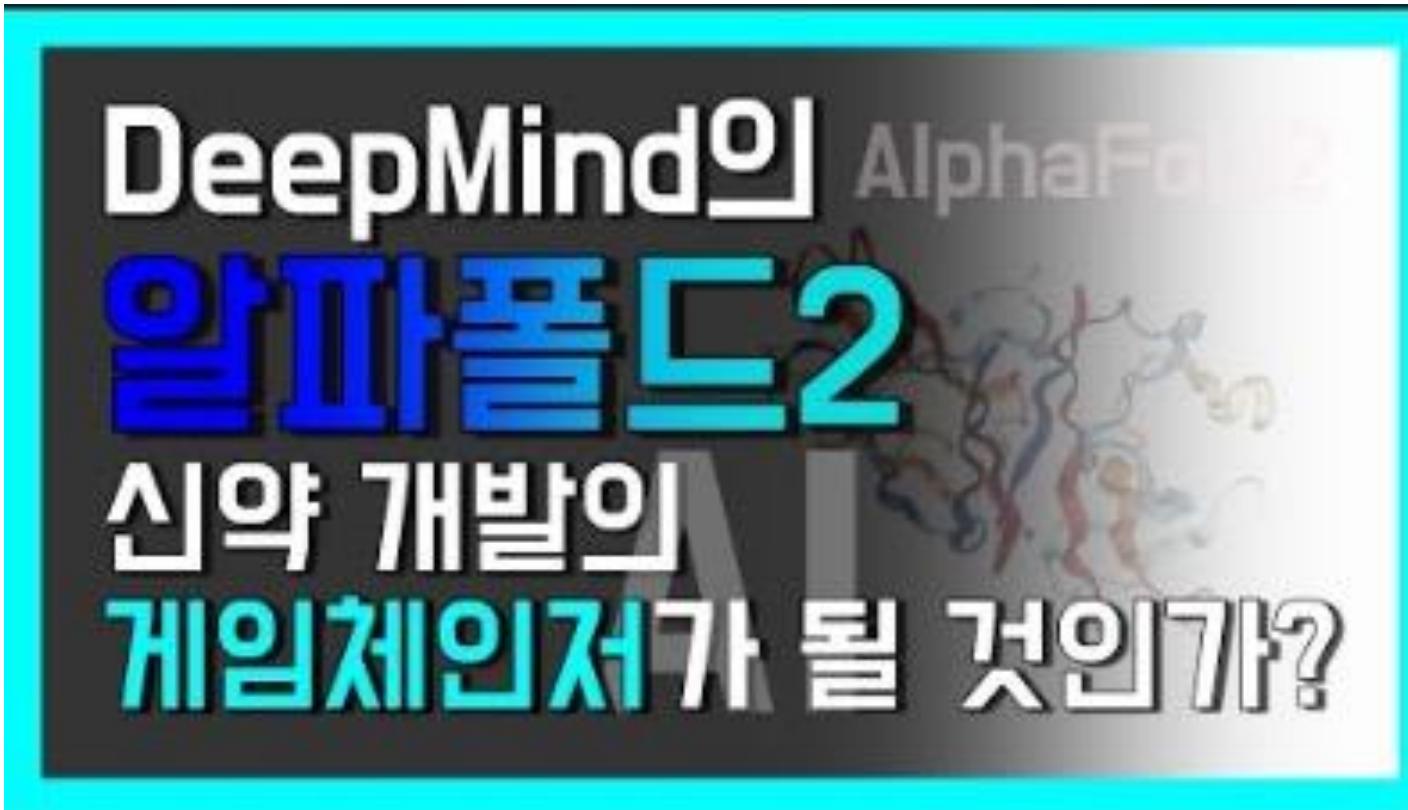
AI Language Modeling : GPT-3

Artificial Intelligence in Context of Human History



Natural Language Processing + Vision : Dall-E

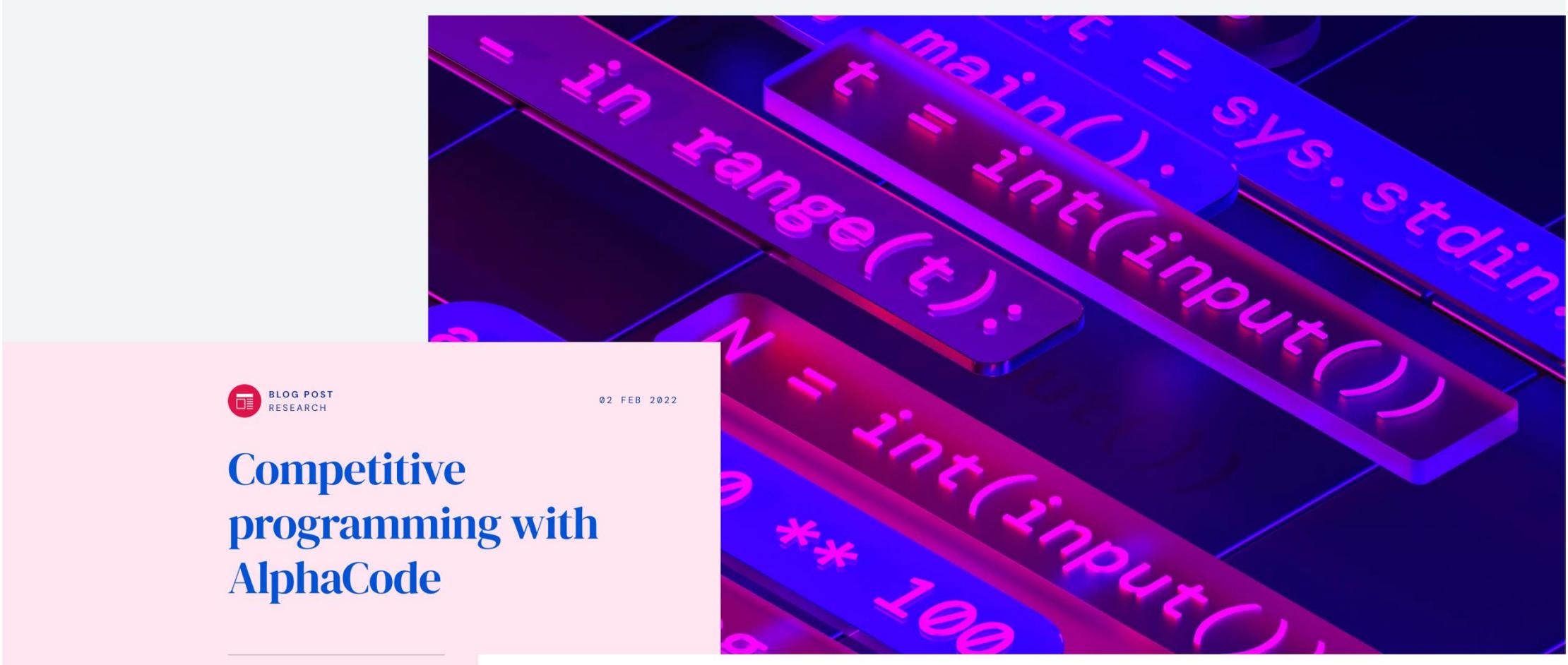
Artificial Intelligence in Context of Human History



Protein Structure Prediction : Alphafold

Artificial Intelligence in Context of Human History

DeepMind > Blog > Competitive programming with AlphaCode



Competitive programming (2022.02.)

Artificial Intelligence in Context of Human History



Dall-E 2 (2022.04.)

Artificial Intelligence in Context of Human History



An art gallery displaying Monet paintings. The art gallery is flooded. Robots are going around the art gallery using paddle boards.



A photo of a Corgi dog riding a bike in Times Square. It is wearing sunglasses and a beach hat.



A strawberry mug filled with white sesame seeds. The mug is floating in a dark chocolate sea.



A robot couple fine dining with Eiffel Tower in the background.

Text-to-Text: Imagen by Google (2022.05.)

Artificial Intelligence in Context of Human History

Stable Diffusion

Stable Diffusion was made possible thanks to a collaboration with Stability AI and Runway and builds upon our previous work:

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach*, Andreas Blattmann*, Dominik Lorenz, Patrick Esser, Björn Ommer

CVPR '22 Oral / GitHub / arXiv / Project page



Stable Diffusion is a latent text-to-image diffusion model. Thanks to a generous compute donation from Stability AI and support from LAION, we were able to train a Latent Diffusion Model on 512x512 images from a subset of the LAION-5B database. Similar to Google's Imagen, this model uses a frozen CLIP ViT-L/14 text encoder to condition the model on text prompts. With its 860M UNet and 123M text encoder, the model is relatively lightweight and runs on a GPU with at least 10GB VRAM. See [this section](#) below and the [model card](#).

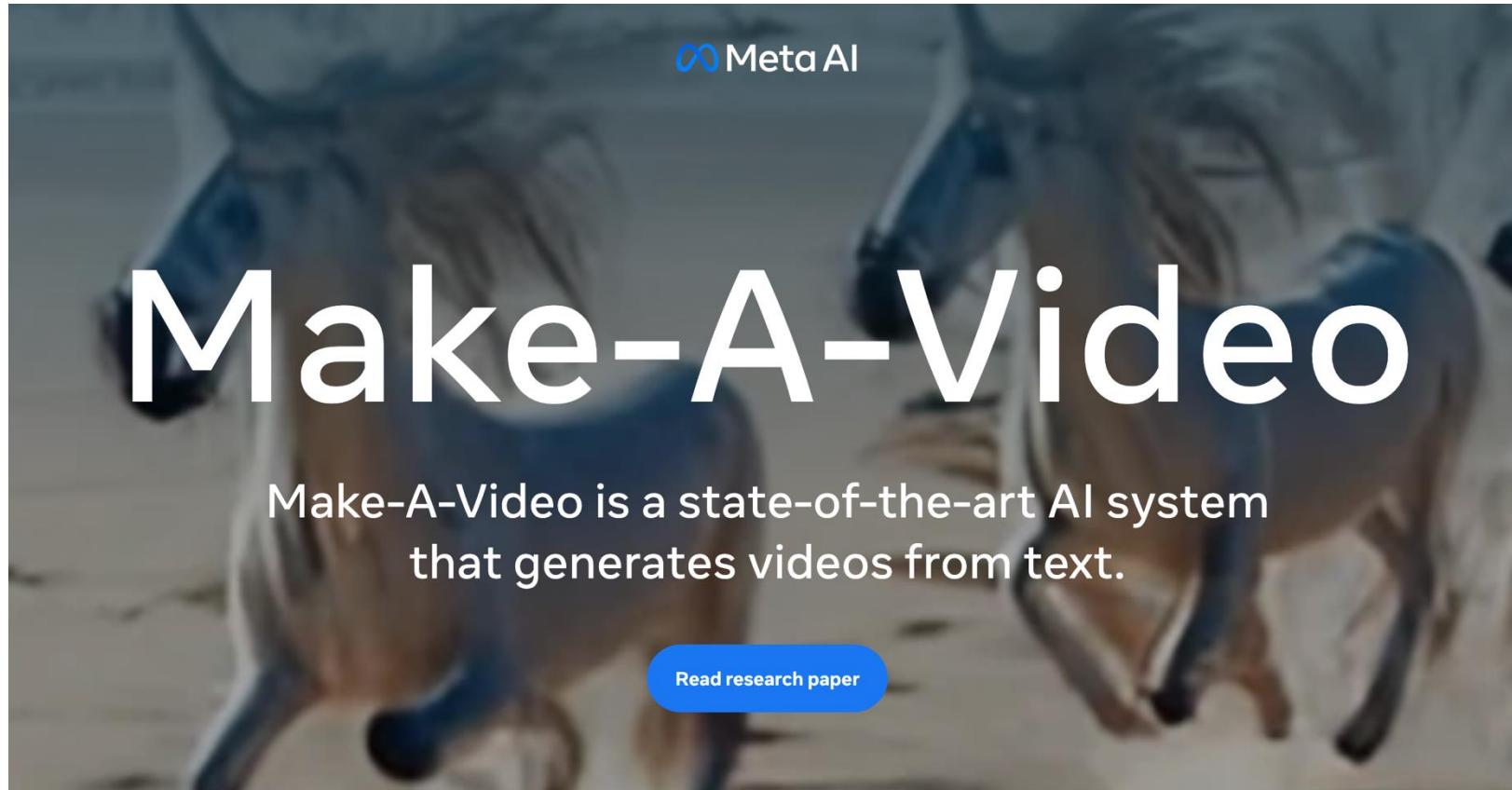
Text-to-Text: Stable Diffusion (2022.08.)

<https://huggingface.co/spaces/stabilityai/stable-diffusion>

<https://stability.ai/blog/stable-diffusion-public-release>

<https://github.com/CompVis/stable-diffusion>

Artificial Intelligence in Context of Human History



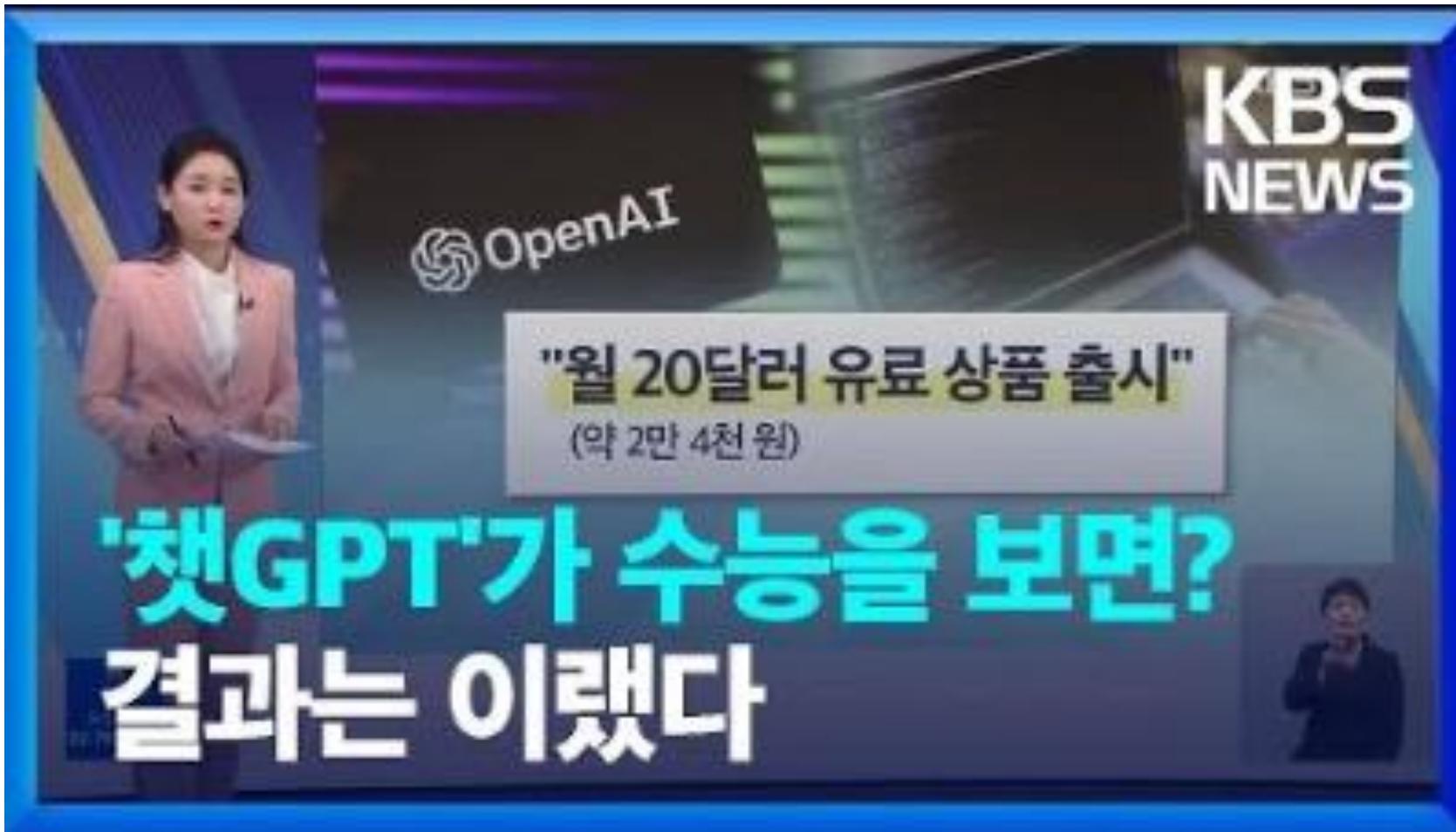
Text-to-Video: Make-A-Video (2022.08.)

Artificial Intelligence in Context of Human History



Text-to-Video: Imagen Video (2022.10.)

Artificial Intelligence in Context of Human History



ChatGPT (2022.11.30.)

Artificial Intelligence in Context of Human History



Text-to-Video Generation (2023.03.)

Artificial Intelligence in Context of Human History



Text-to-Video Generation (2023.03.)

Artificial Intelligence in Context of Human History

ANTHROPIC

Claude ▾ Research Company Careers News

Meet Claude

Claude is AI for all of us. Whether you're brainstorming alone or building with a team of thousands, Claude is here to help.

Try Claude

Get API access



Claude by Anthropic (2023.04.)

Artificial Intelligence in Context of Human History



Stable Diffusion XL (2023.07)

Artificial Intelligence in Context of Human History



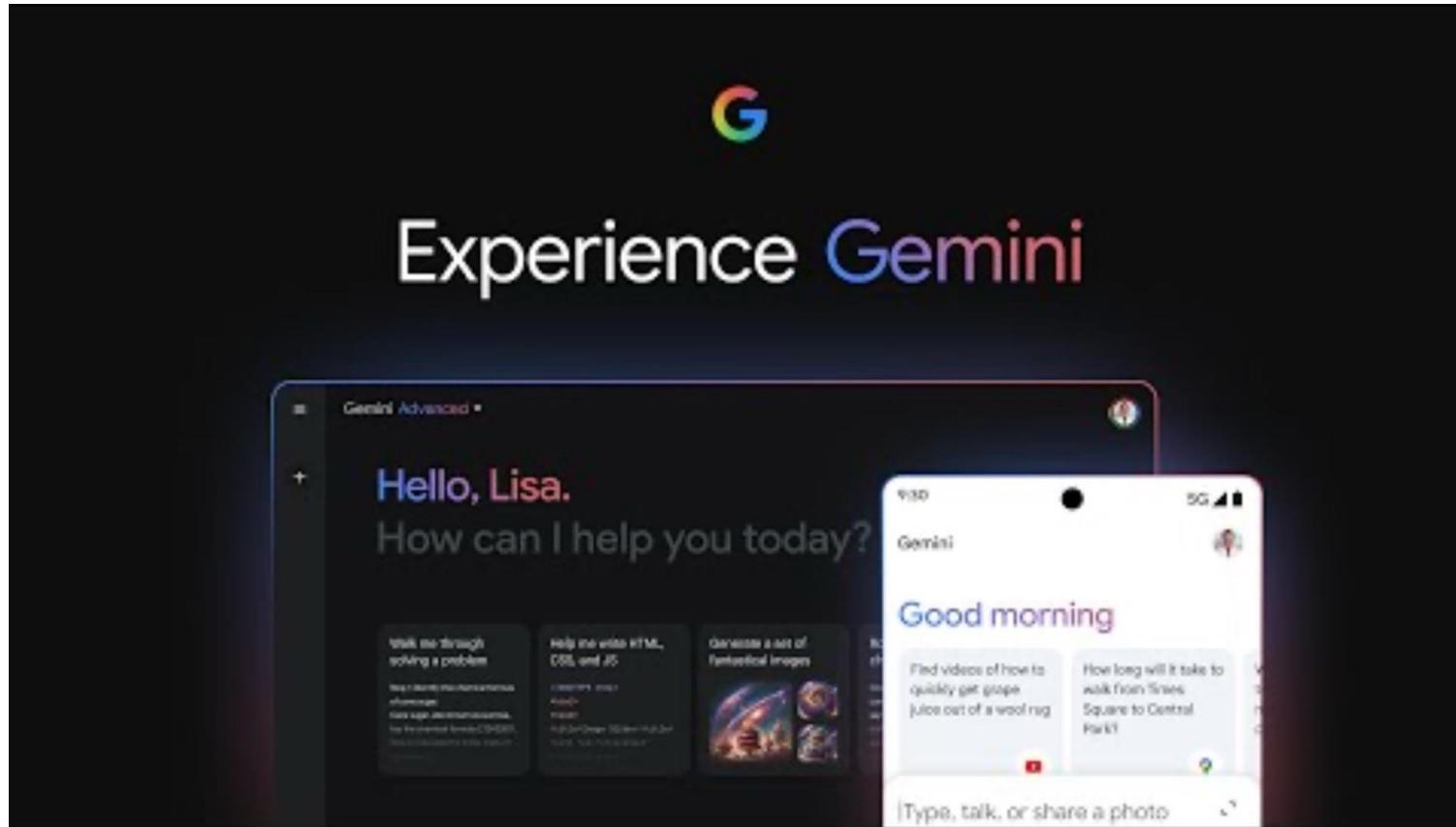
Stable Diffusion Video (2023.11)

Artificial Intelligence in Context of Human History



DALLE-3 (2023.12)

Artificial Intelligence in Context of Human History



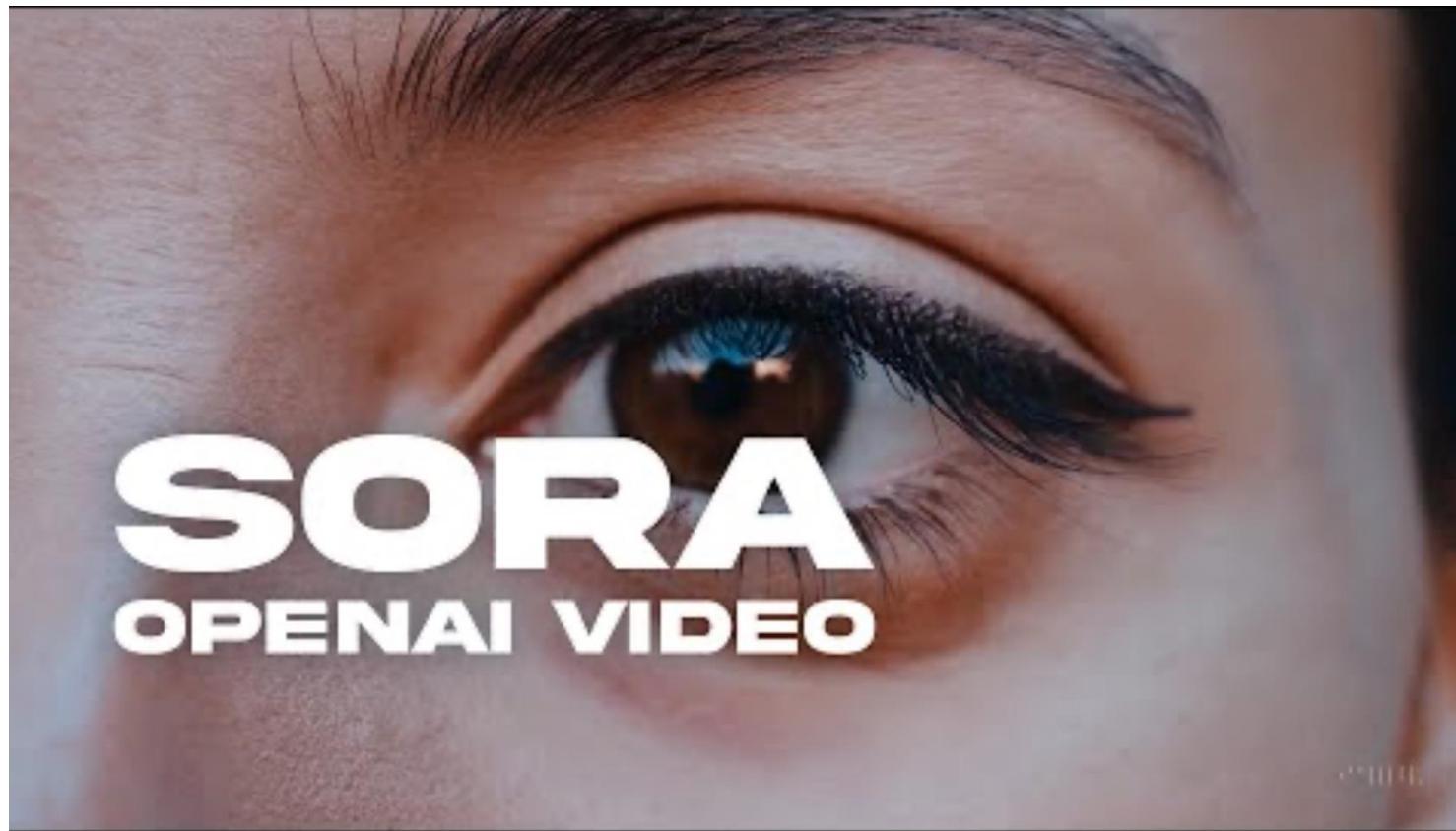
Gemini (2023.12)

Artificial Intelligence in Context of Human History



Lumiere (2024.01.23)

Artificial Intelligence in Context of Human History



SORA (2024.02.16)

Artificial Intelligence in Context of Human History



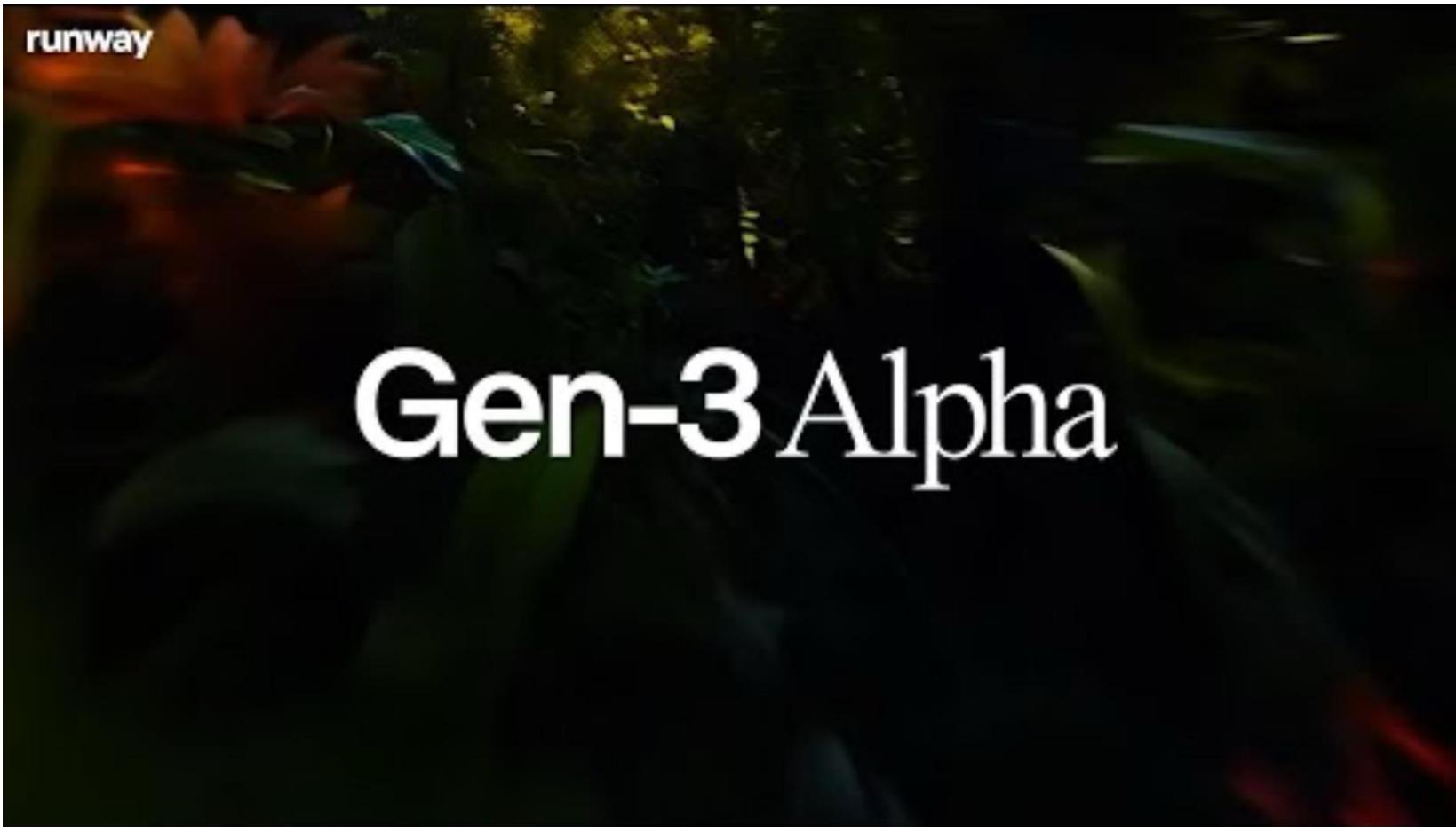
Stable Diffusion 3 (2024.02.22)

Artificial Intelligence in Context of Human History



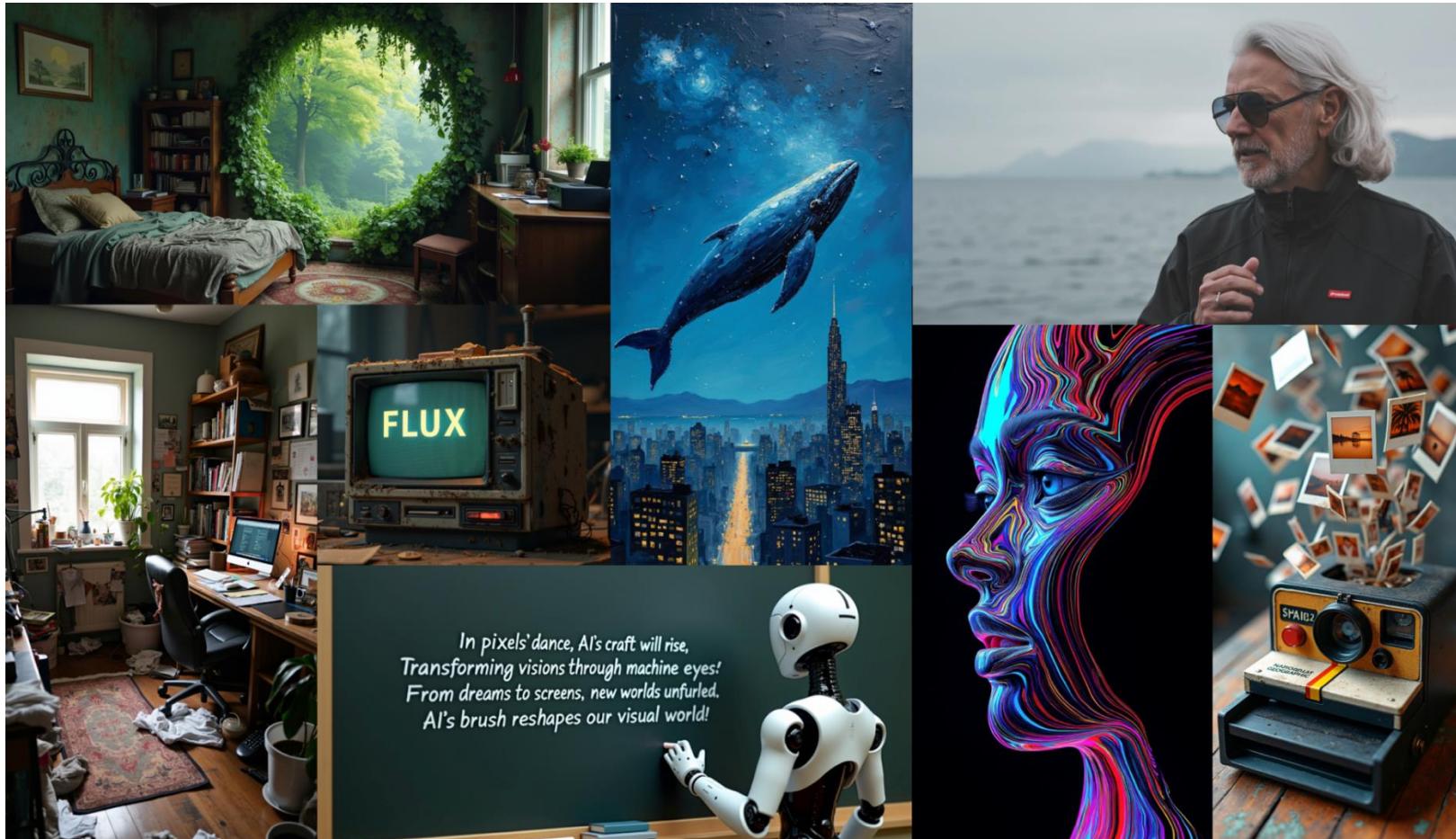
GPT-4o (2024.05.13)

Artificial Intelligence in Context of Human History



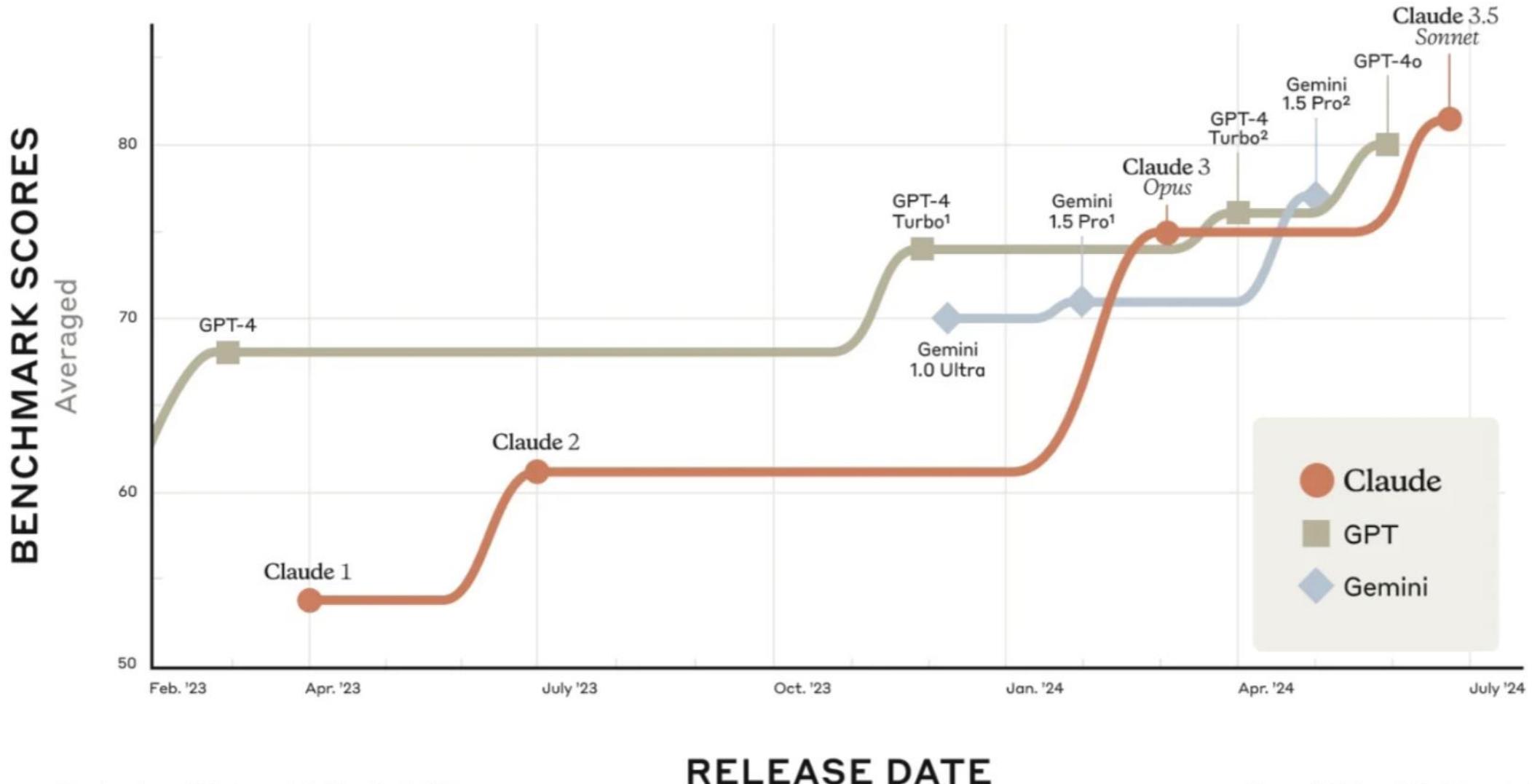
GEN-3 (2024.06.17)

Artificial Intelligence in Context of Human History



FLUX (2024.08.01)

AI model release and capabilities timeline



Averaged benchmarks are highest reported without best-of-N:
MMLU, GPQA, MATH, MGSM, DROP F1, HumanEval pass@1,
MMMU, AI2D, ChartQA, DocQA, Mathvista

Source: Publicly available data; evaluation scores are
the average of representative scores found online.
1 = Initial release; 2 = Second release

History of Deep Learning Ideas and Milestones

- 1943: Neural networks
- 1957-62: Perceptron
- 1970-86: Backpropagation, RBM, RNN
- 1979-98: CNN, MNIST, LSTM, Bidirectional RNN
- 2006: “Deep Learning”, DBN
- 2009: ImageNet
- 2012: AlexNet
- 2014: GANs
- 2016: ResNet
- 2016-17: AlphaGo, AlphaZero
- 2017: Transformers
- 2020: Diffusion model
- 2021: Vision Transformer
- 2022: chatGPT

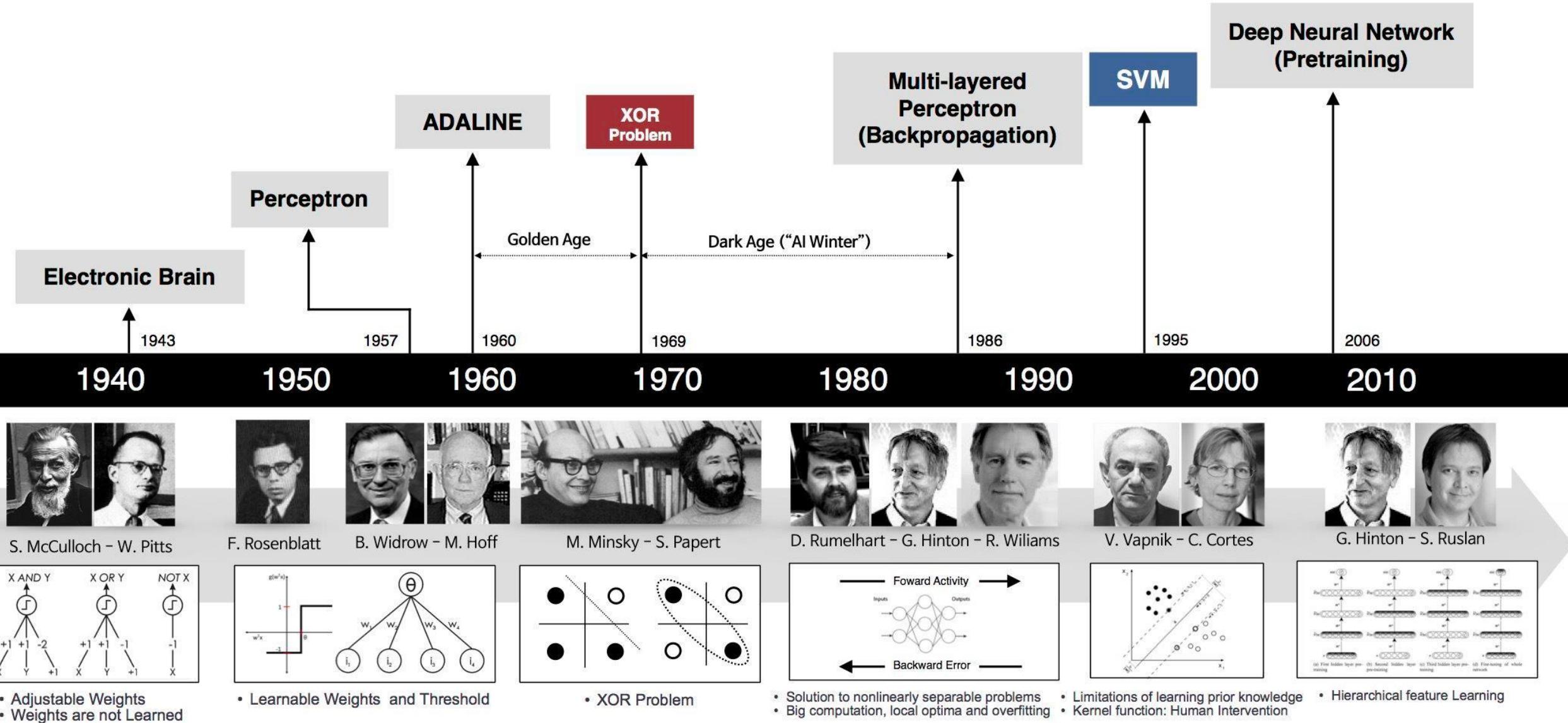
* Dates are for perspective and not as definitive historical record of invention or credit

History of Deep Learning Ideas and Milestones

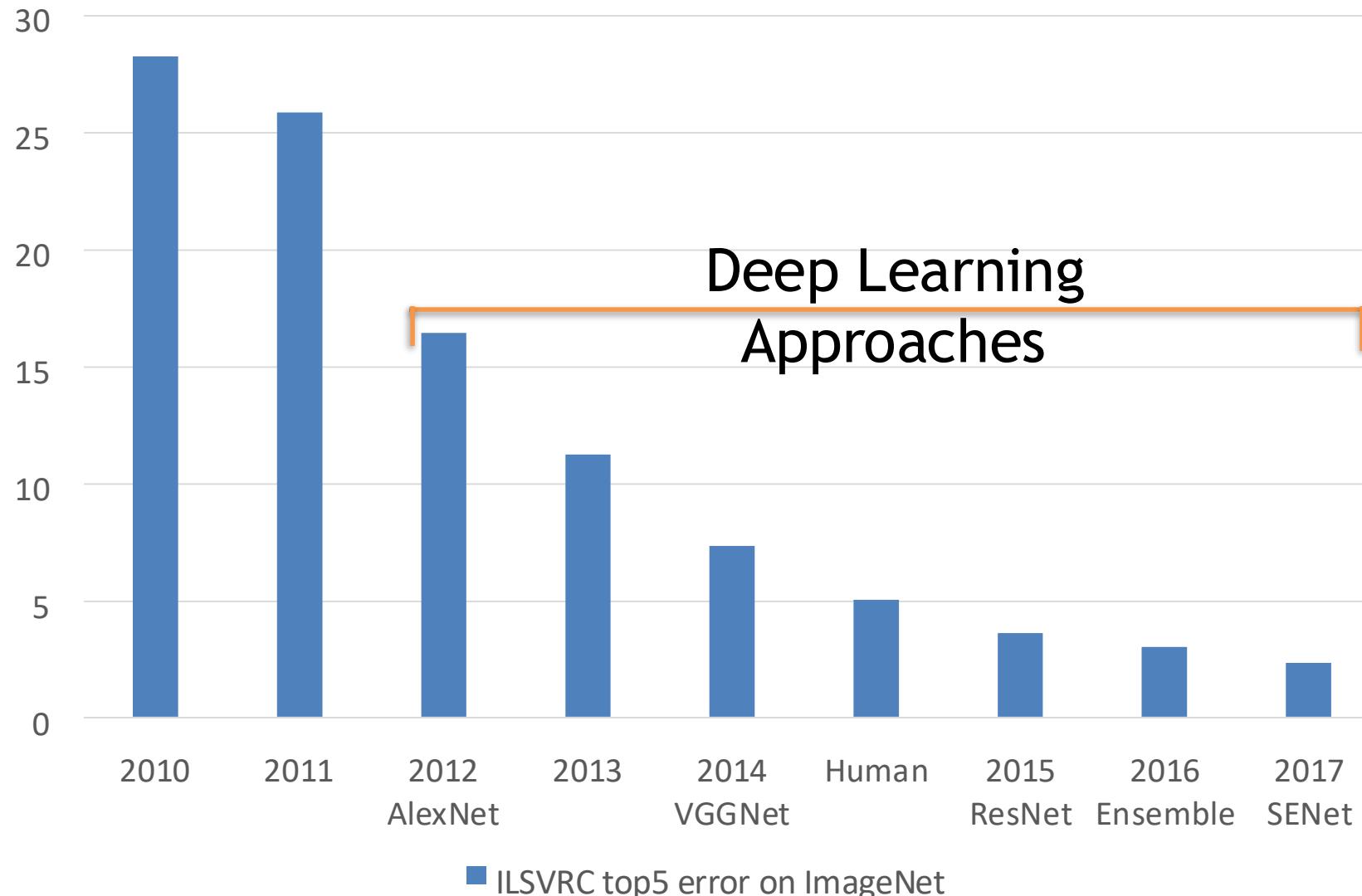
- 1943: Neural networks
- 1957-62: Perceptron
- 1970-86: Backpropagation, RBM, RNN
- 1979-98: CNN, MNIST, LSTM, Bidirectional RNN
- 2006: “Deep Learning”, DBN
- 2009: ImageNet
- 2012: **AlexNet**
- 2014: GANs
- 2016: **ResNet**
- 2016-17: AlphaGo, AlphaZero
- 2017: **Transformers**
- 2020: **Diffusion model**
- 2021: **Vision Transformer**
- 2022: **chatGPT**

* Dates are for perspective and not as definitive historical record of invention or credit

History of Deep Learning Ideas and Milestones

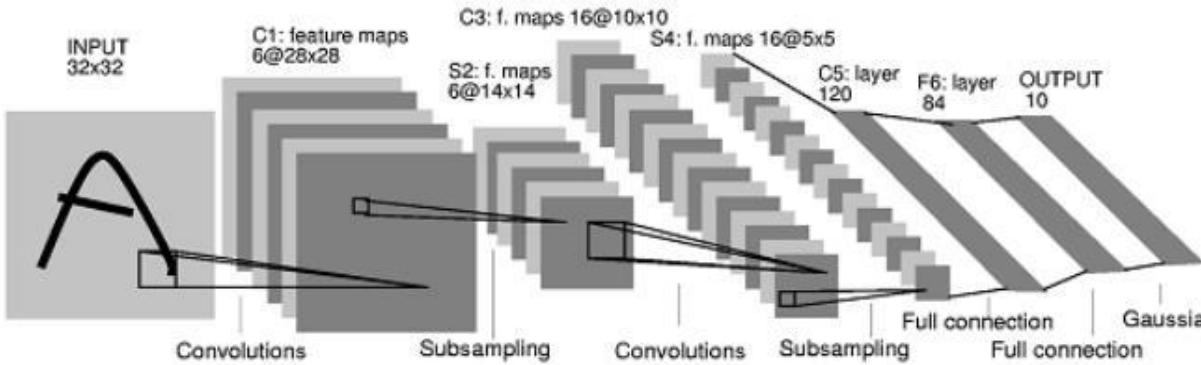


Breakthrough



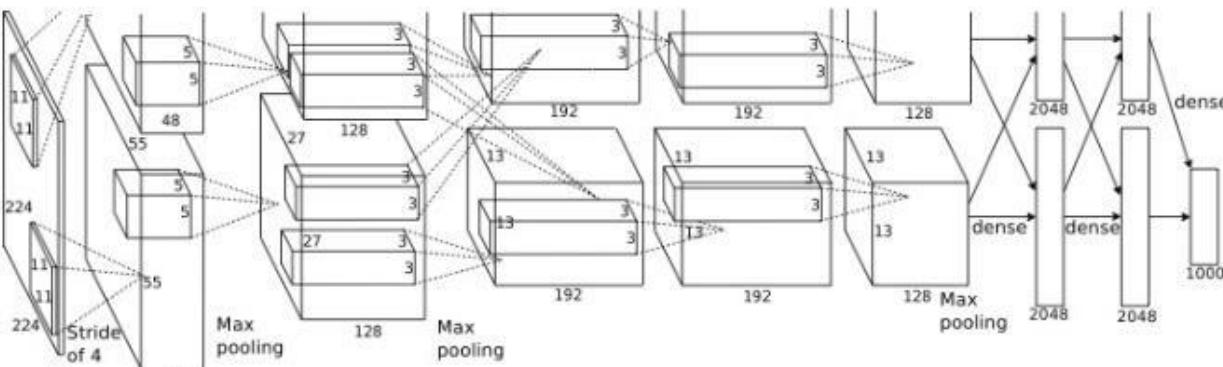
What has changed?

1998
LeCun
et al.



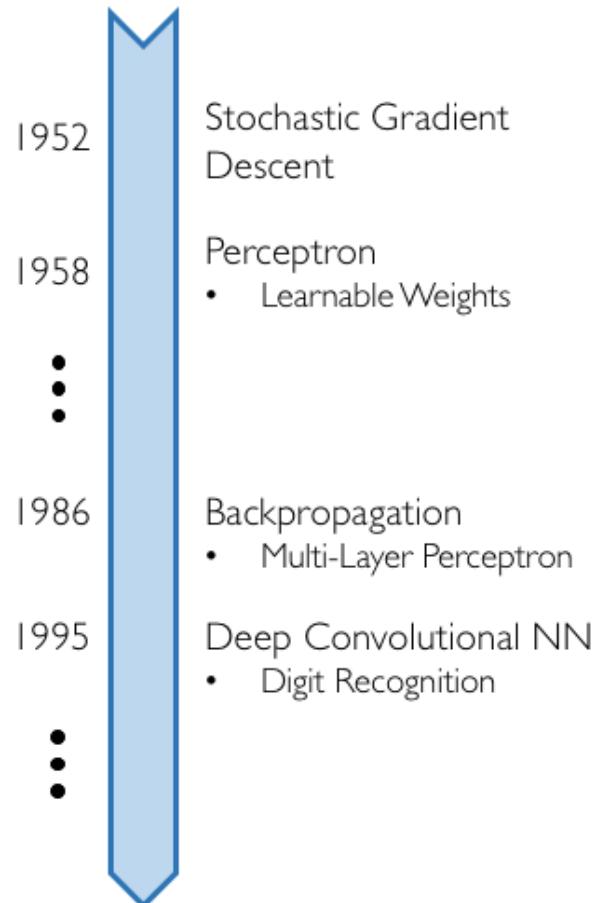
- MNIST digit recognition dataset
- 10^7 pixels used in training

2012
Krizhevsky
et al.



- ImageNet image recognition dataset
- 10^{14} pixels used in training

Why Now?



Neural Networks date back decades, so why the resurgence?

I. Big Data

- Larger Datasets
- Easier Collection & Storage



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



3. Software

- Improved Techniques
- New Models
- Toolboxes



Turing Award for Deep Learning



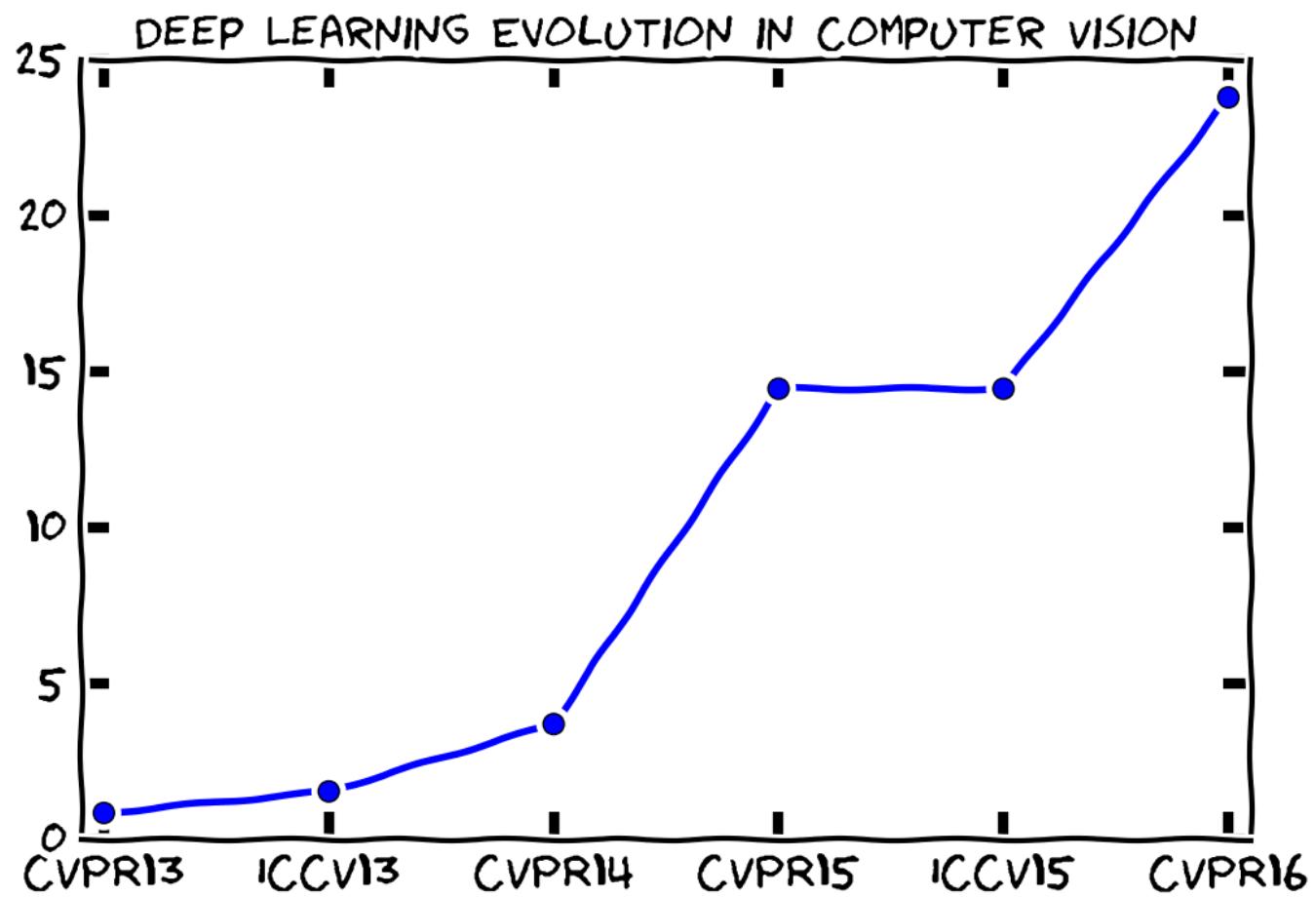
- Yann LeCun
- Geoffrey Hinton
- Yoshua Bengio

Turing Award given for:

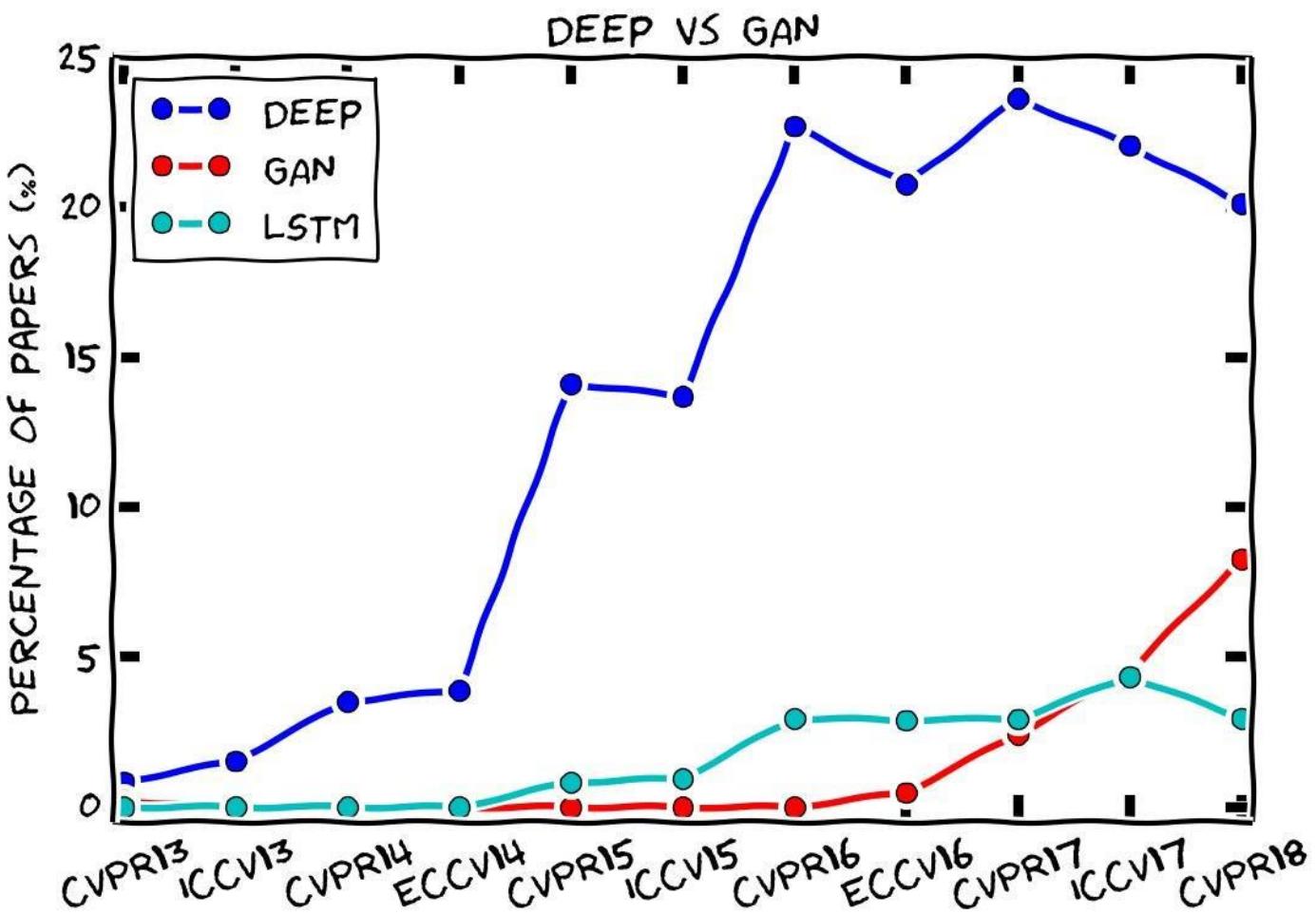
- “The conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing.”
- (Also, for popularization in the face of skepticism.)

Google Scholar Rank (h-index)

Categories	Publication	<u>h5-index</u>	<u>h5-median</u>
1.	Nature	<u>488</u>	745
2.	IEEE/CVF Conference on Computer Vision and Pattern Recognition CVPR	<u>440</u>	689
3.	The New England Journal of Medicine	<u>434</u>	897
4.	Science	<u>409</u>	633
5.	Nature Communications	<u>375</u>	492
6.	The Lancet	<u>368</u>	678
7.	Neural Information Processing Systems NeurIPS	<u>337</u>	614
8.	Advanced Materials	<u>327</u>	420
9.	Cell	<u>320</u>	482
10.	International Conference on Learning Representations ICLR	<u>304</u>	584



Credits: Dr. Pont-Tuset, ETH Zurich



Credits: Dr. Pont-Tuset, ETH Zurich

What is Deep Learning?

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



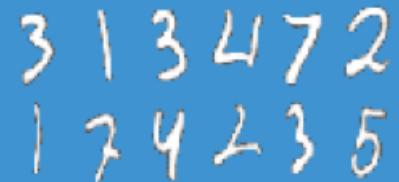
MACHINE LEARNING

Ability to learn without explicitly being programmed

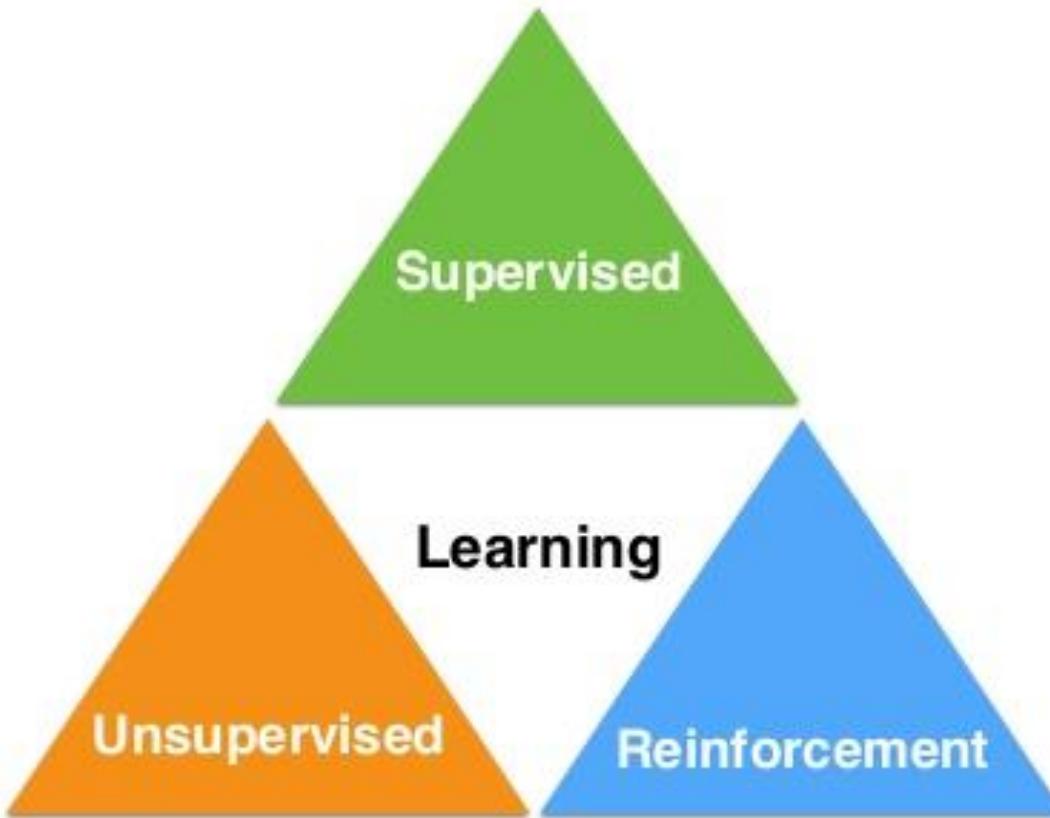


DEEP LEARNING

Extract patterns from data using neural networks



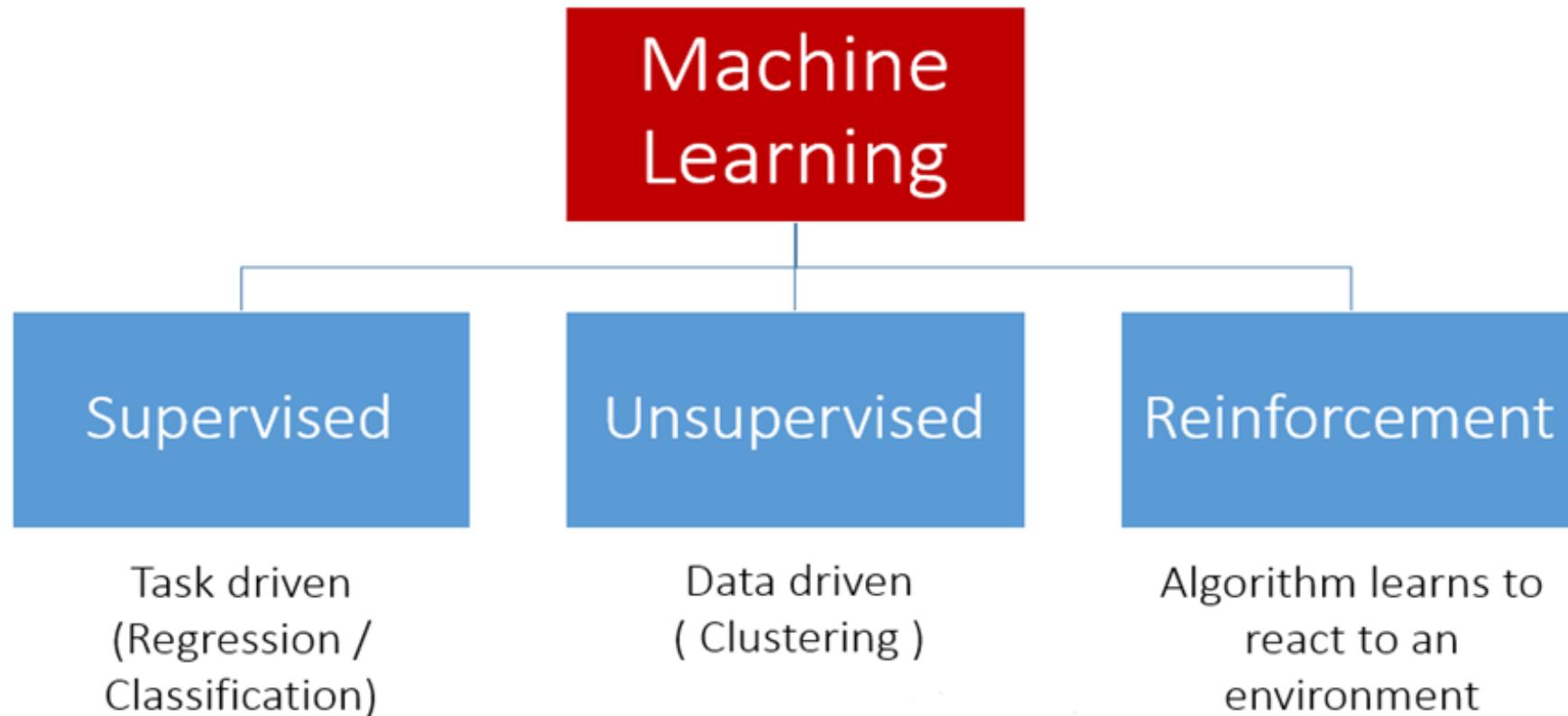
- Labeled data
- Direct feedback
- Predict outcome/future



- No labels
- No feedback
- “Find hidden structure”

- Decision process
- Reward system
- Learn series of actions

Types of Machine Learning



Types of Machine Learning

Supervised Learning



Classification

- Fraud detection
- Email Spam Detection
- Diagnostics
- Image Classification

Unsupervised Learning



Dimensionality Reduction

- Text Mining
- Face Recognition
- Big Data Visualization
- Image Recognition

Reinforcement Learning



- Gaming
- Finance Sector
- Manufacturing
- Inventory Management
- Robot Navigation

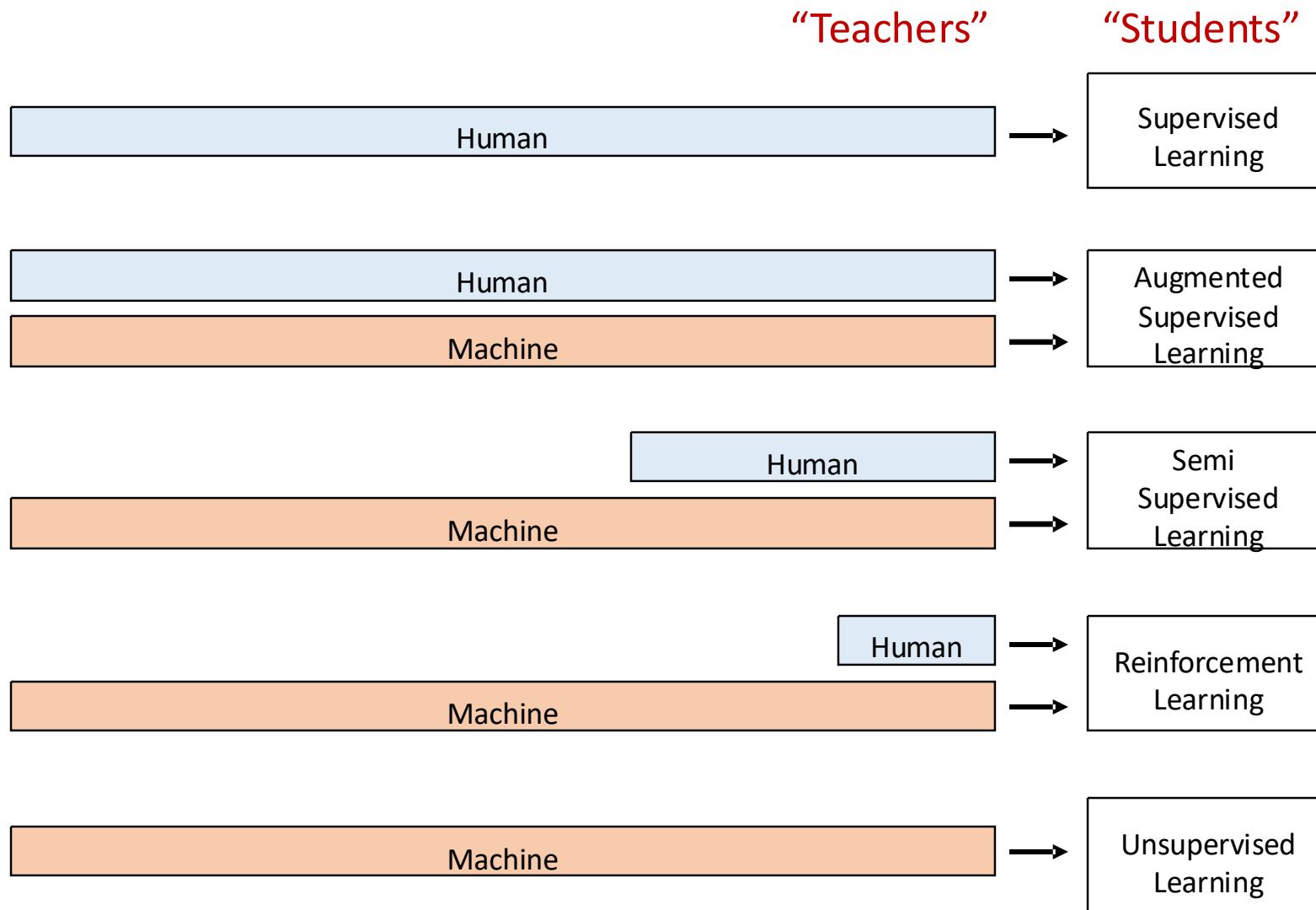
Regression

- Risk Assessment
- Score Prediction

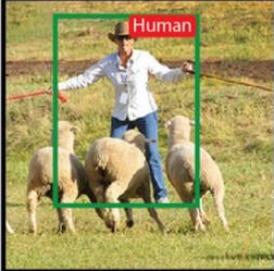
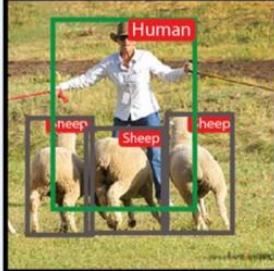
Clustering

- Biology
- City Planning
- Targetted Marketing

Deep Learning from Human and Machine



Supervised Learning

	Image Classification Classify an image based on the dominant object inside it. datasets: MNIST, CIFAR, ImageNet		Object Localization Predict the image region that contains the dominant object. Then image classification can be used to recognize object in the region datasets: ImageNet		Object Recognition Localize and classify all objects appearing in the image. This task typically includes: proposing regions then classify the object inside them. datasets: PASCAL, COCO
	Semantic Segmentation Label each pixel of an image by the object class that it belongs to, such as human, sheep, and grass in the example. datasets: PASCAL, COCO		Instance Segmentation Label each pixel of an image by the object class and object instance that it belongs to. datasets: PASCAL, COCO		Keypoint Detection Detect locations of a set of predefined keypoints of an object, such as keypoints in a human body, or a human face. datasets: COCO

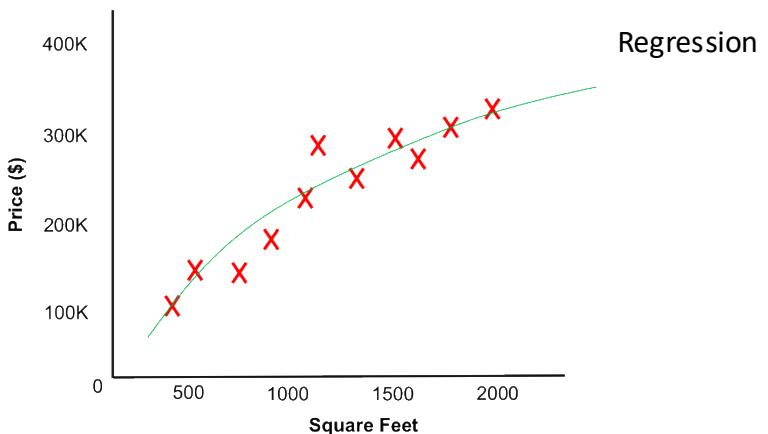


Image captioning

"man in black shirt is playing guitar."

Unsupervised Learning



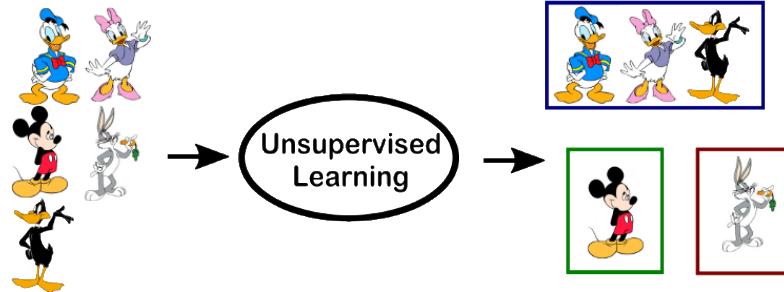
GAN



Image-to-image translation



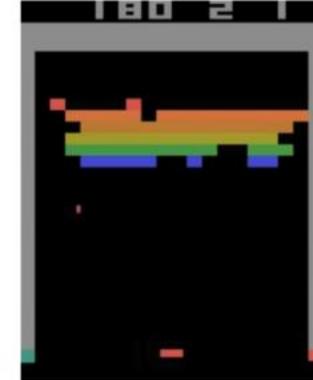
Clustering



Reinforcement Learning



(a) 슈퍼마리오



(b) Atari - 벽돌깨기

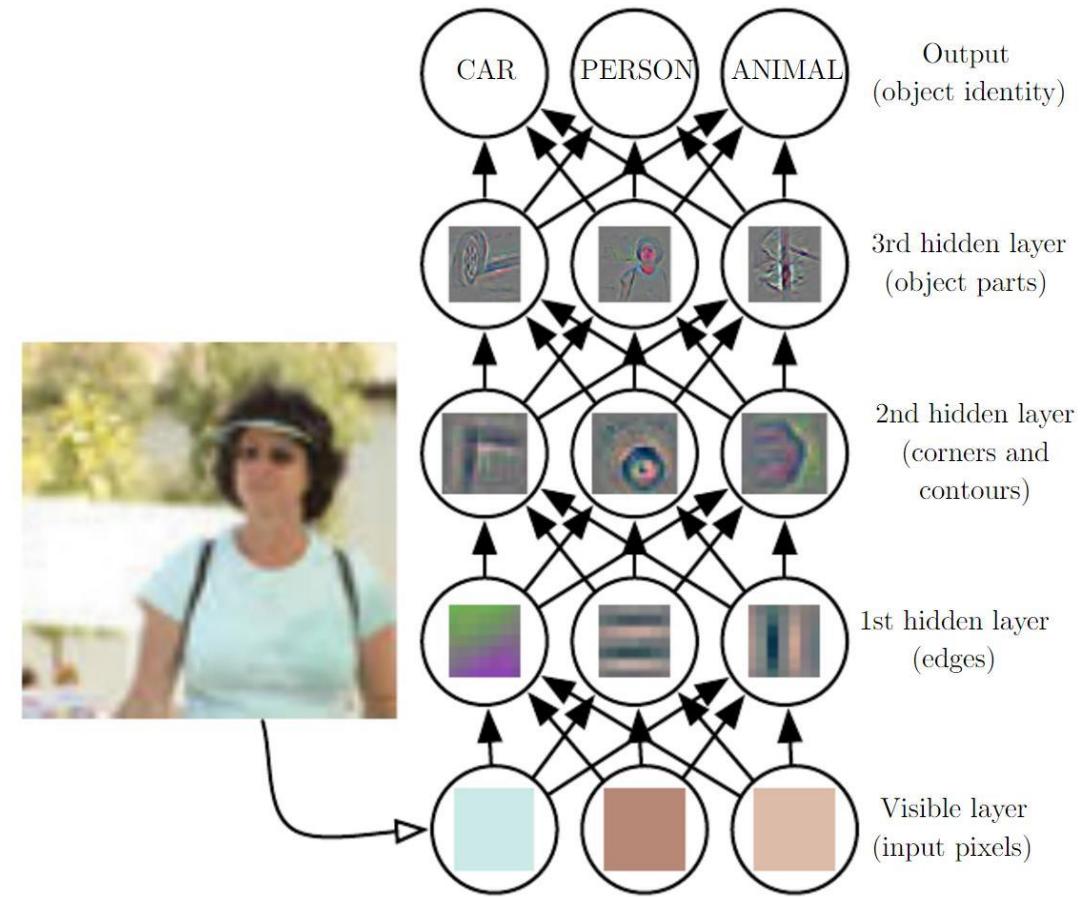
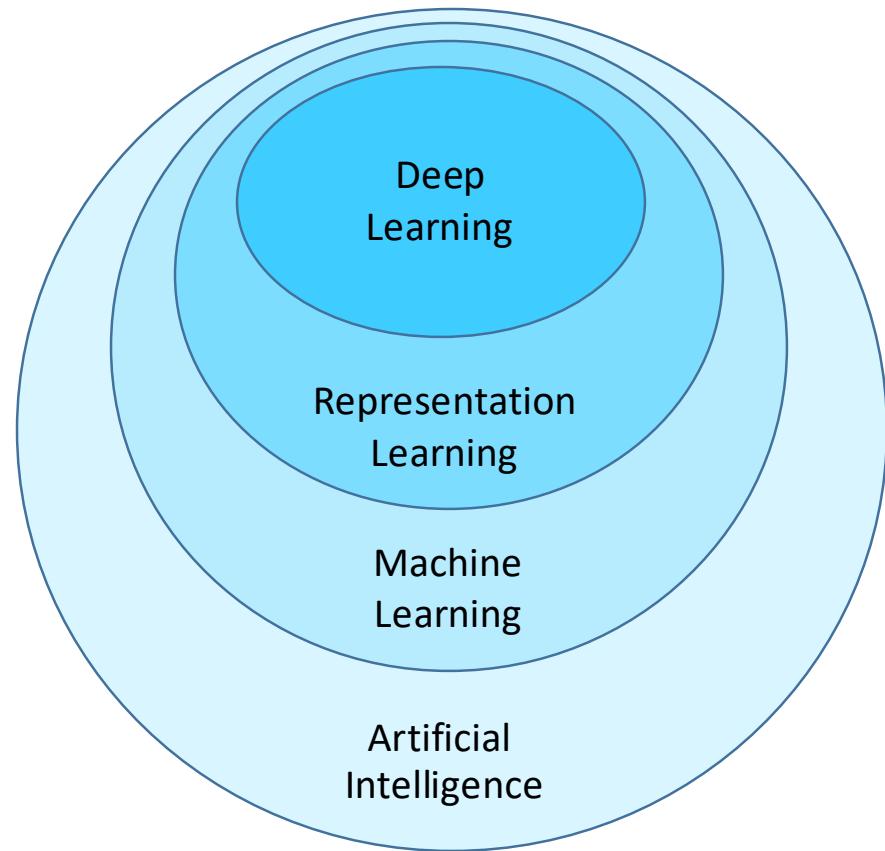


스타크래프트2
프로게이머

vs

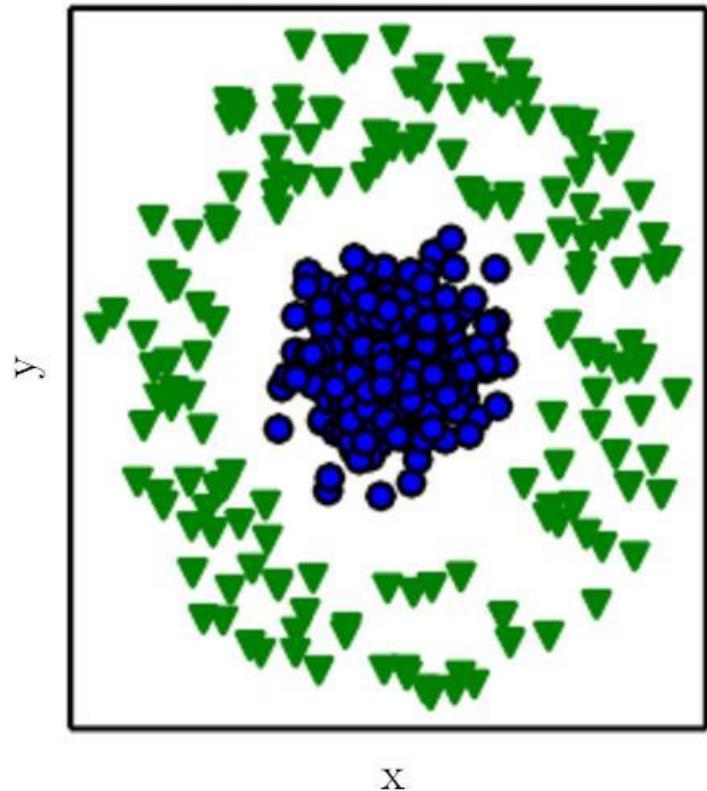
DeepMind

Deep Learning is Representation Learning (aka Feature Learning)

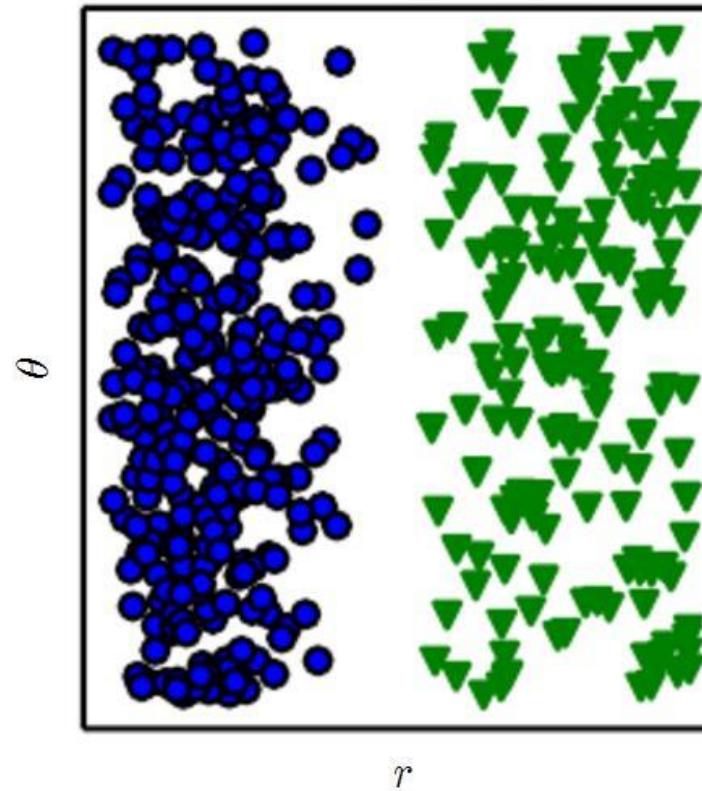


Representation Matters

Cartesian coordinates

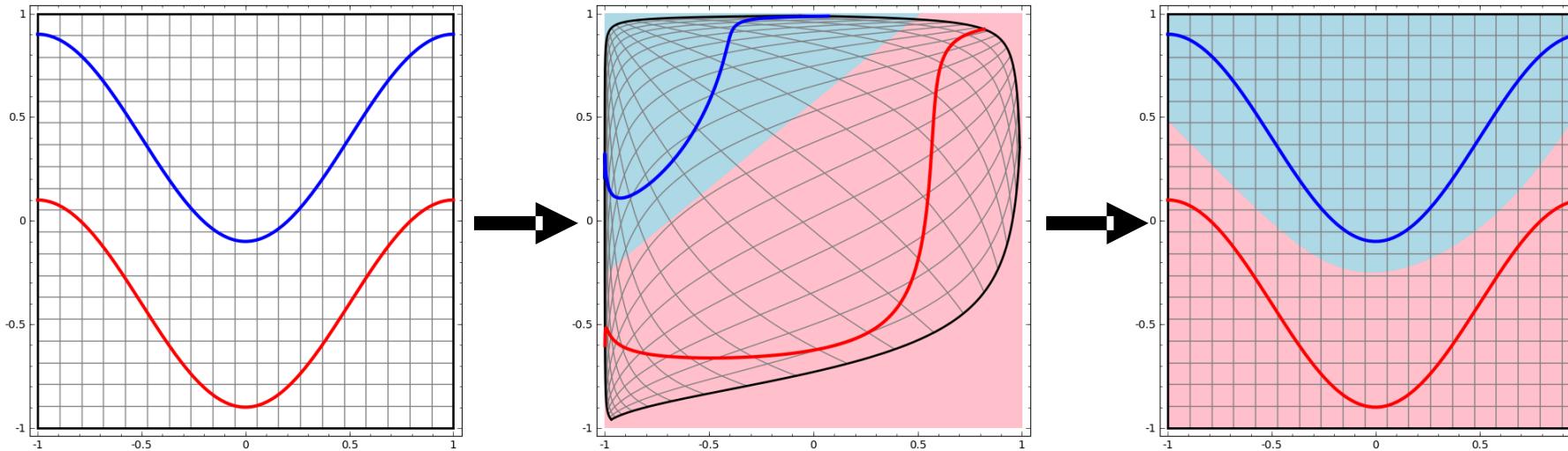


Polar coordinates



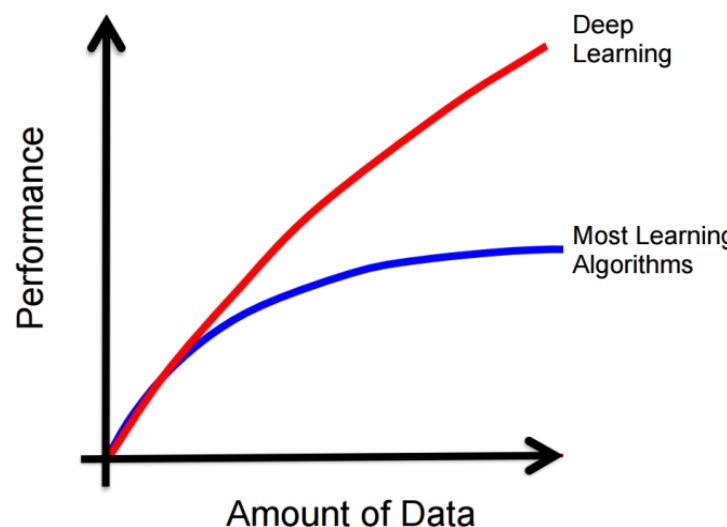
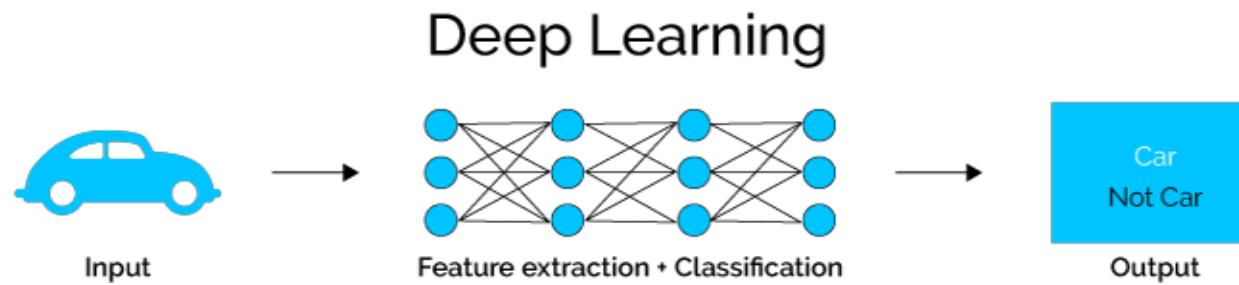
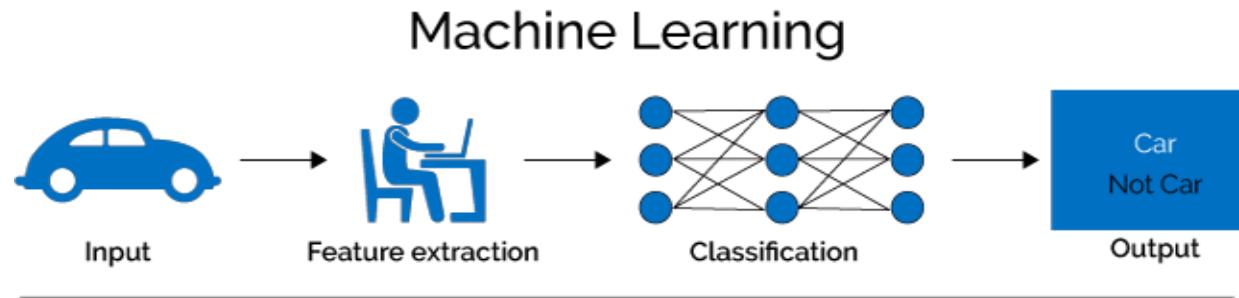
Task: Draw a line to separate the green triangles and blue circles.

Deep Learning is Representation Learning (aka Feature Learning)



Task: Draw a line to separate the **blue curve** and **red curve**

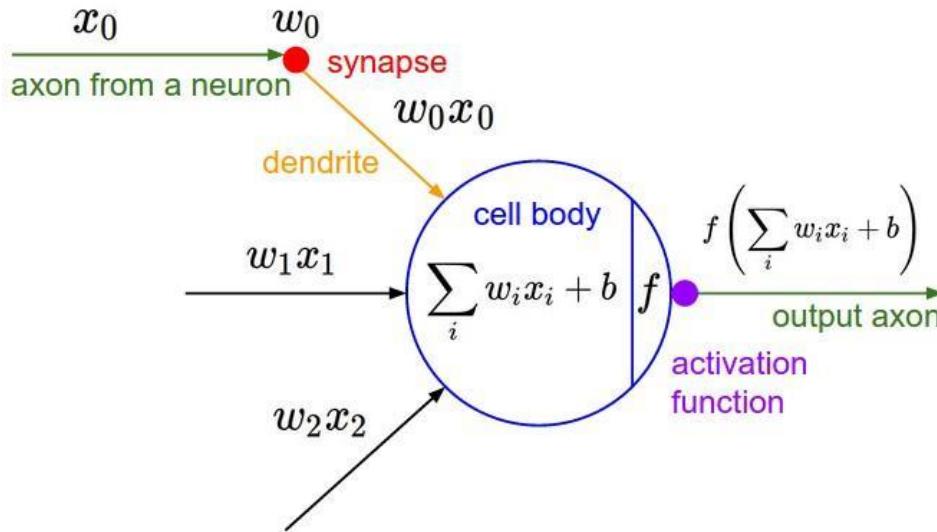
Why Deep Learning? Scalable Machine Learning



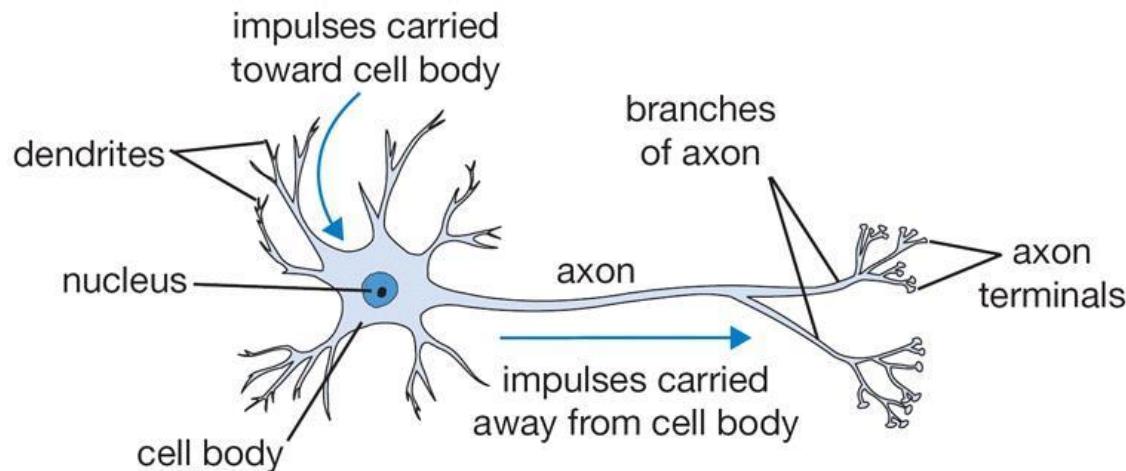
What is Neural Network?

From perceptron to structural build block of deep learning

Neuron: Biological Inspiration for Computation



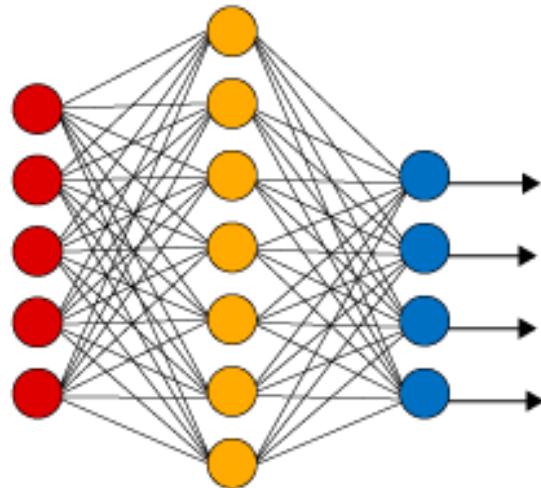
(Artificial) Neuron:
computational building block for
the “neural network”



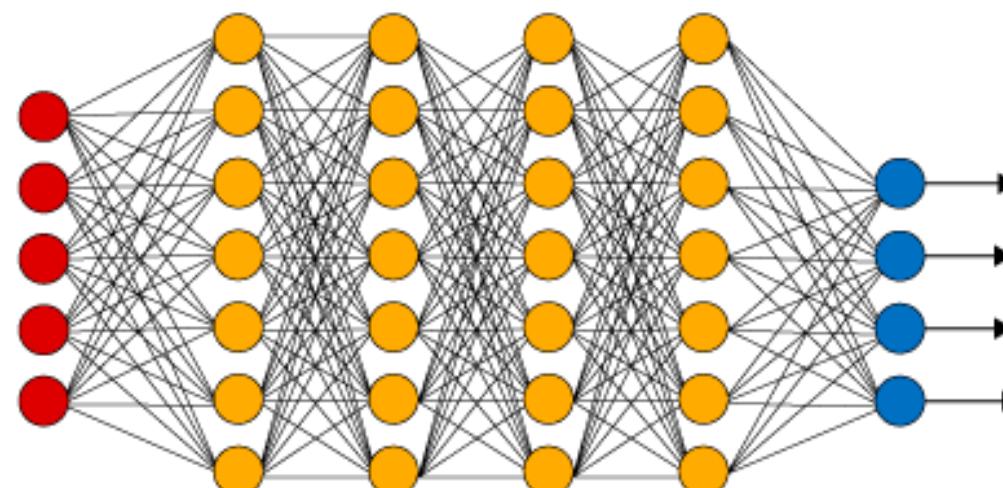
Neuron:
computational building block for
the brain

Combing Neurons in Hidden Layers: The “Emergent” Power to Approximate

Simple Neural Network



Deep Learning Neural Network



● Input Layer

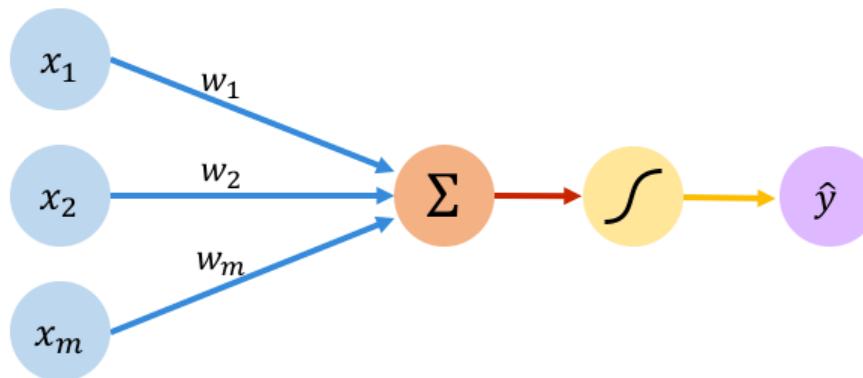
● Hidden Layer

● Output Layer

Universal Approximation Theorem:

For any arbitrary function $f(x)$,
there exists a neural network that closely approximate it for any input x

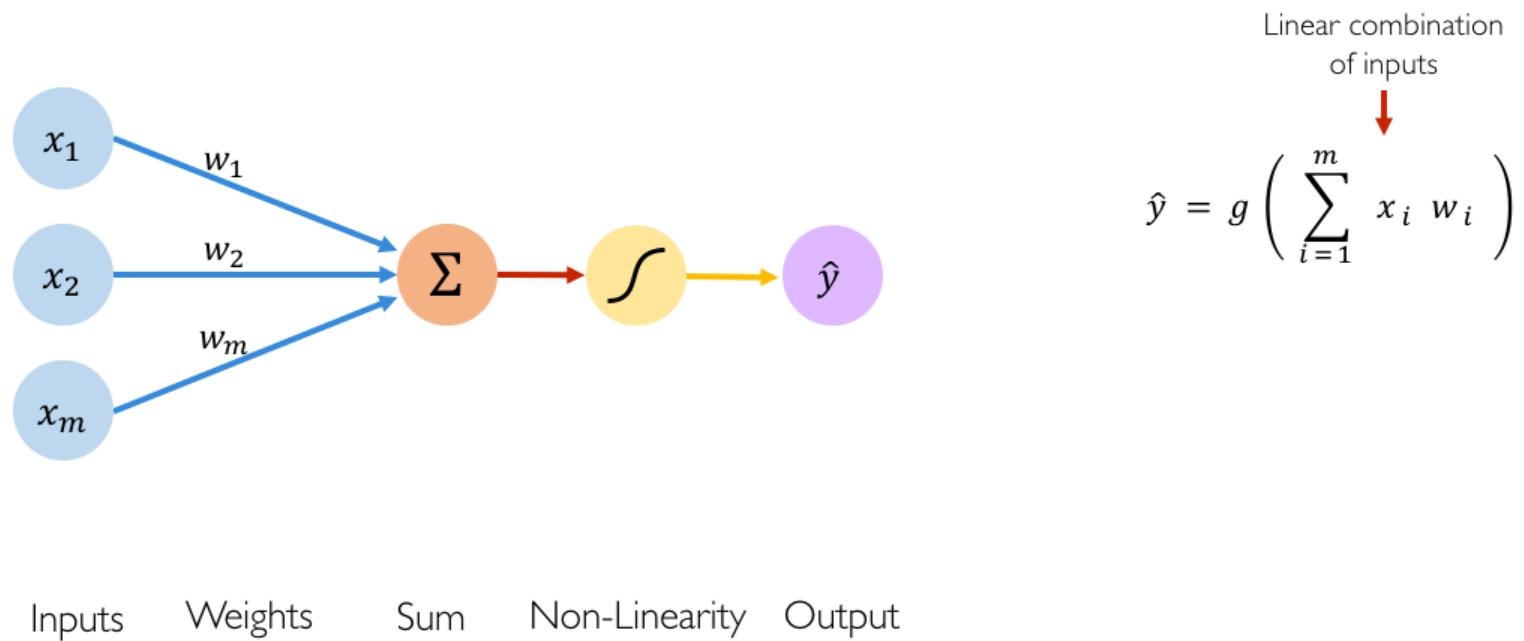
The Perceptron: Forward Propagation



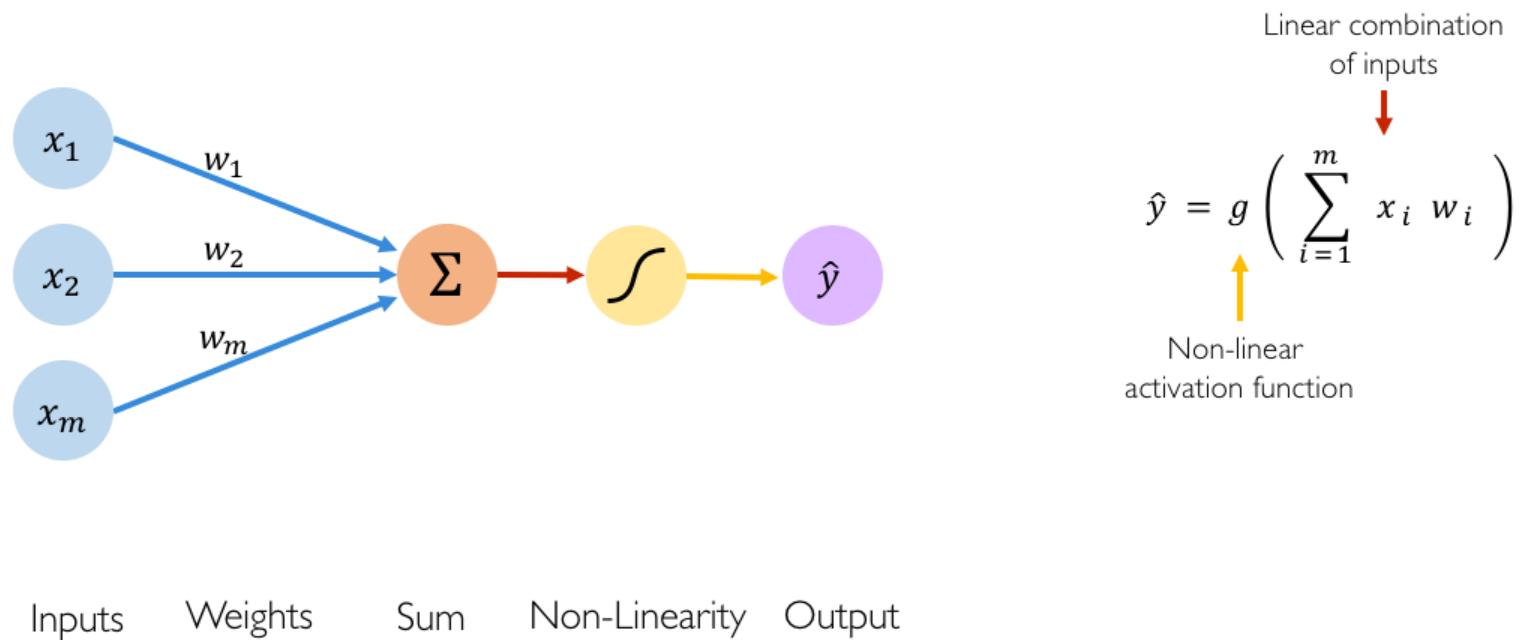
$$\hat{y} = g \left(\sum_{i=1}^m x_i w_i \right)$$

Inputs Weights Sum Non-Linearity Output

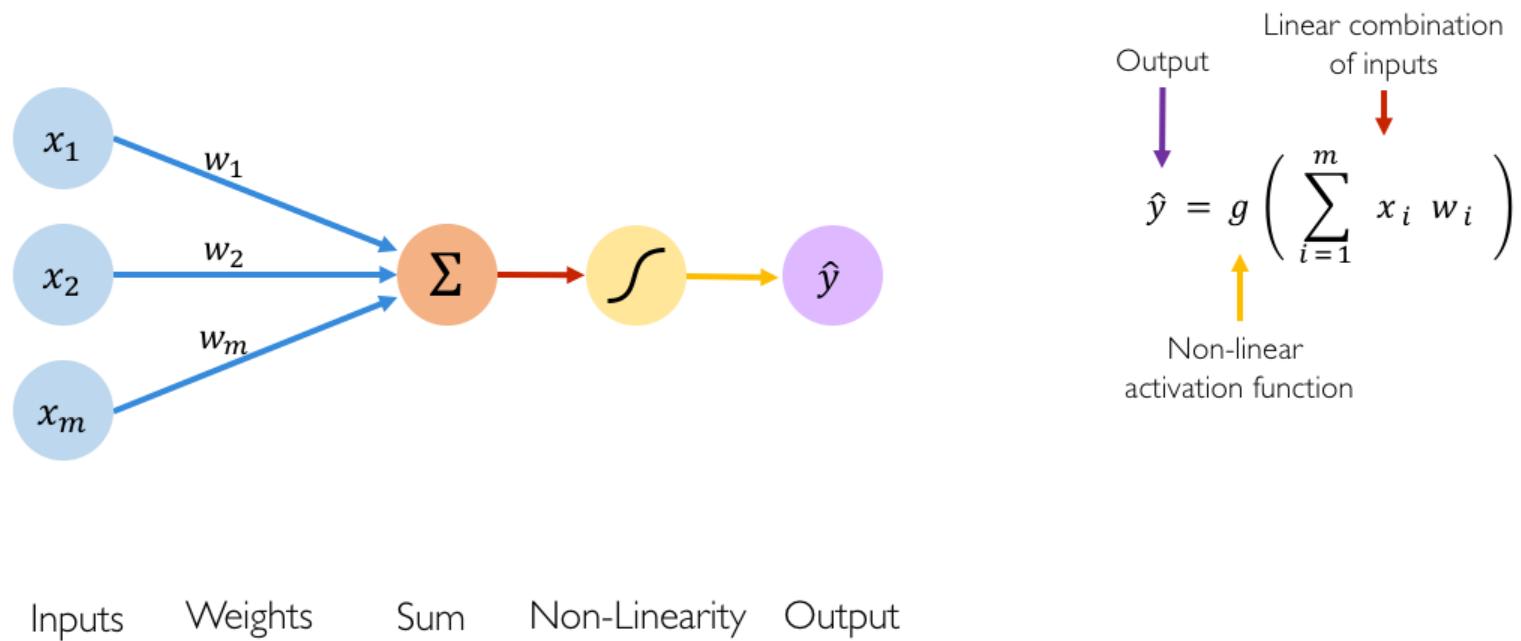
The Perceptron: Forward Propagation



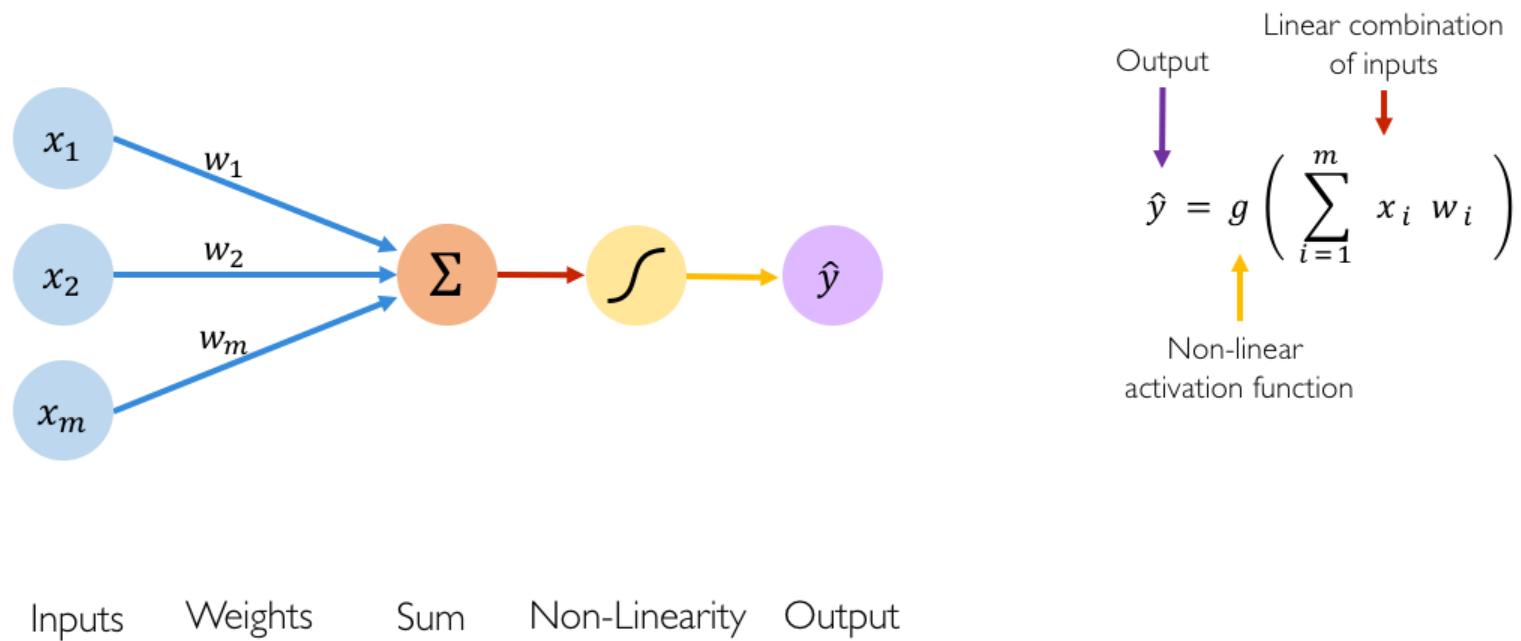
The Perceptron: Forward Propagation



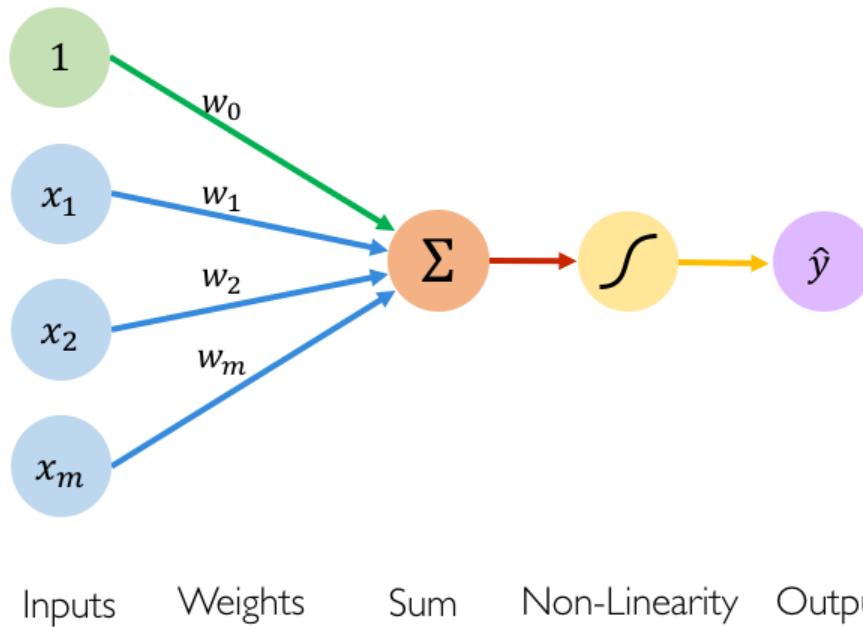
The Perceptron: Forward Propagation



The Perceptron: Forward Propagation



The Perceptron: Forward Propagation



Linear combination of inputs

Output

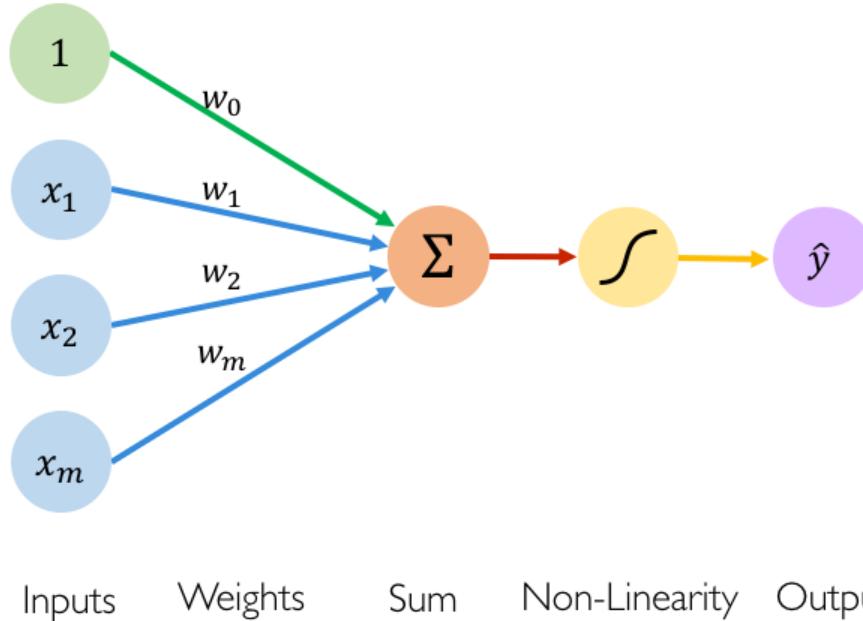
$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$

Non-linear activation function

Bias

Diagram illustrating the mathematical formula for the perceptron's output. The output \hat{y} is the result of applying the activation function g to the linear combination of inputs and weights. The linear combination is the bias w_0 plus the sum of the products of each input x_i and its corresponding weight w_i .

The Perceptron: Forward Propagation

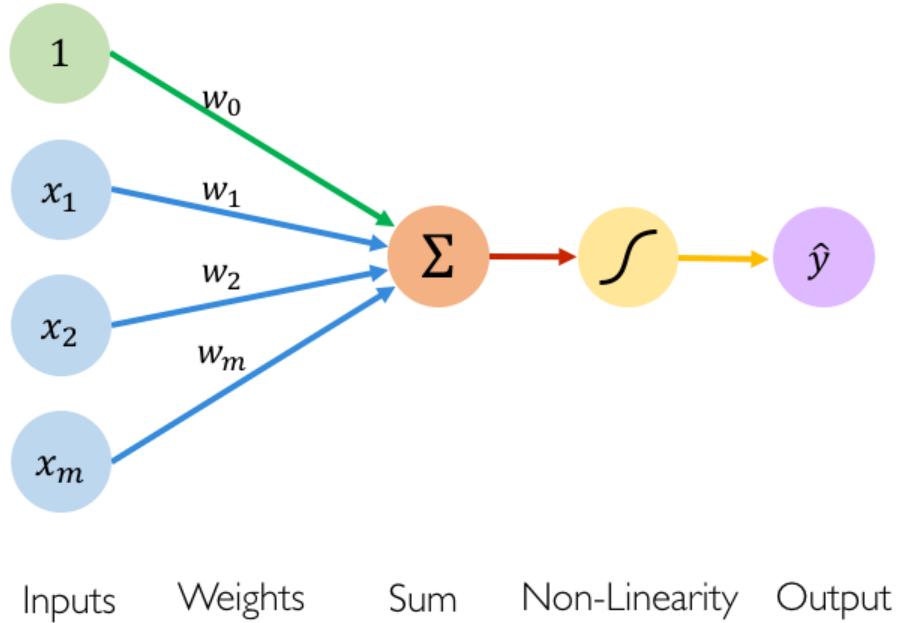


$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

$$\hat{y} = g (w_0 + \mathbf{X}^T \mathbf{W})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

The Perceptron: Forward Propagation

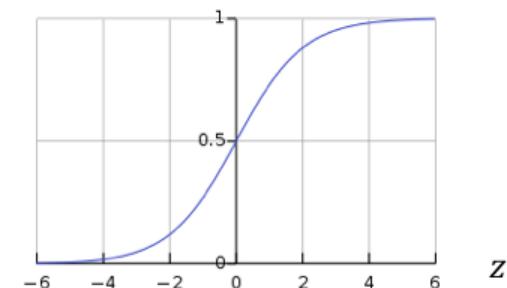


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

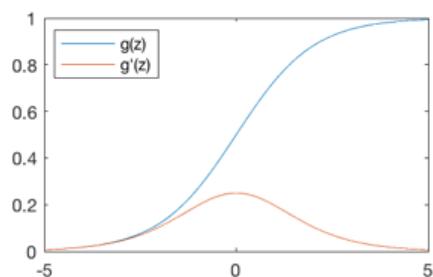
- Example: sigmoid function

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Common Activation Functions

Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

SIGMOID

CLASS `torch.nn.Sigmoid`

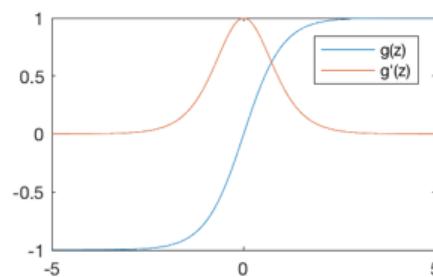
Applies the element-wise function:

`torch.nn.functional.sigmoid(input) → Tensor`

Applies the element-wise function $\text{Sigmoid}(x) = \frac{1}{1+\exp(-x)}$

See [Sigmoid](#) for more details.

Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

TANH

CLASS `torch.nn.Tanh`

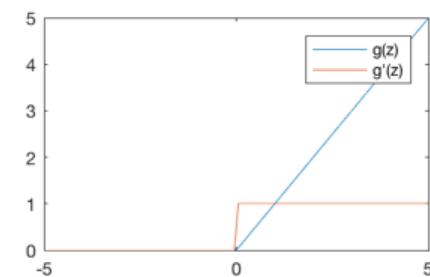
Applies the element-wise function:

`torch.nn.functional.tanh(input) → Tensor`

Applies element-wise, $\text{Tanh}(x) = \tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$

See [Tanh](#) for more details.

Rectified Linear Unit (ReLU)



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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

RELU

CLASS `torch.nn.ReLU(inplace: bool = False)`

Applies the rectified linear unit function element-wise:

`torch.nn.functional.relu(input, inplace=False) → Tensor`

Applies the rectified linear unit function element-wise. See [ReLU](#) for more details.

Dance Moves of Deep Learning Activation Functions

Sigmoid



$$y = \frac{1}{1+e^{-x}}$$

Tanh



$$y = \tanh(x)$$

Step Function



$$y = \begin{cases} 0, & x < n \\ 1, & x \geq n \end{cases}$$

Softplus



$$y = \ln(1+e^x)$$

ReLU



$$y = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

Softsign



$$y = \frac{x}{(1+|x|)}$$

ELU



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

Log of Sigmoid



$$y = \ln\left(\frac{1}{1+e^{-x}}\right)$$

Swish



$$y = \frac{x}{1+e^{-x}}$$

Sinc



$$y = \frac{\sin(x)}{x}$$

Leaky ReLU



$$y = \max(0.1x, x)$$

Mish

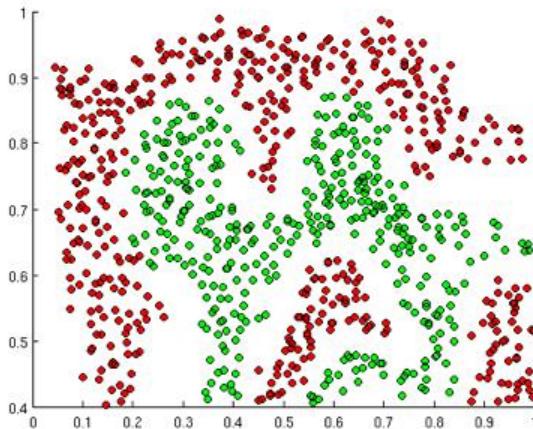


$$y = x(\tanh(\text{softplus}(x)))$$

source: sefiks

Importance of Activation Functions

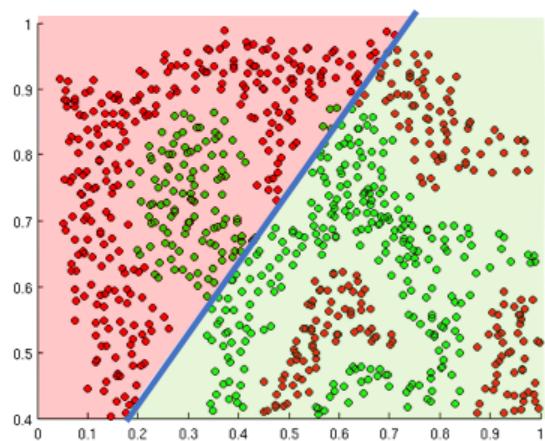
*The purpose of activation functions is to **introduce non-linearities** into the network*



What if we wanted to build a neural network to
distinguish green vs red points?

Importance of Activation Functions

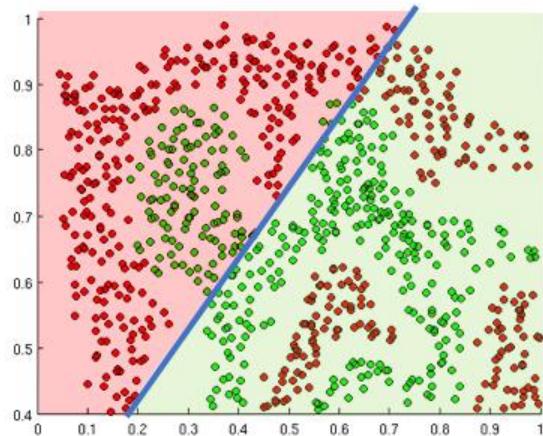
The purpose of activation functions is to **introduce non-linearities** into the network



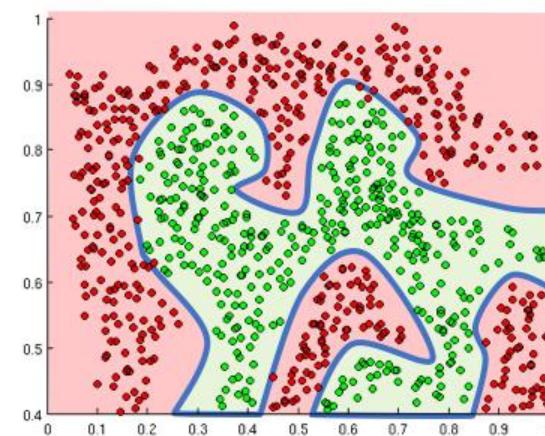
Linear activation functions produce linear decisions no matter the network size

Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network

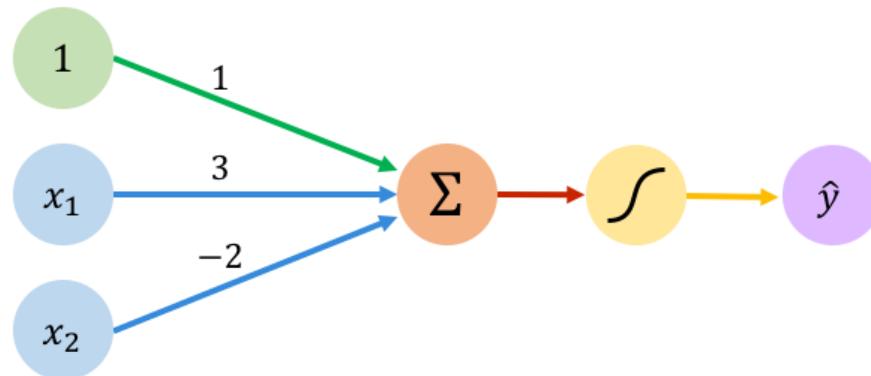


Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

The perceptron : Example

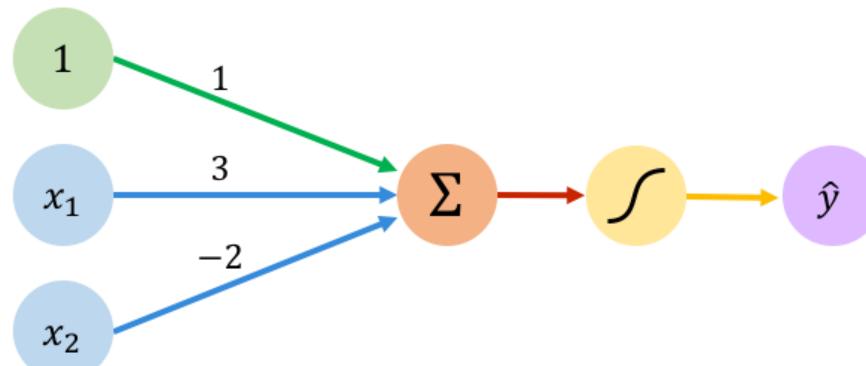


We have: $w_0 = 1$ and $\mathbf{w} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

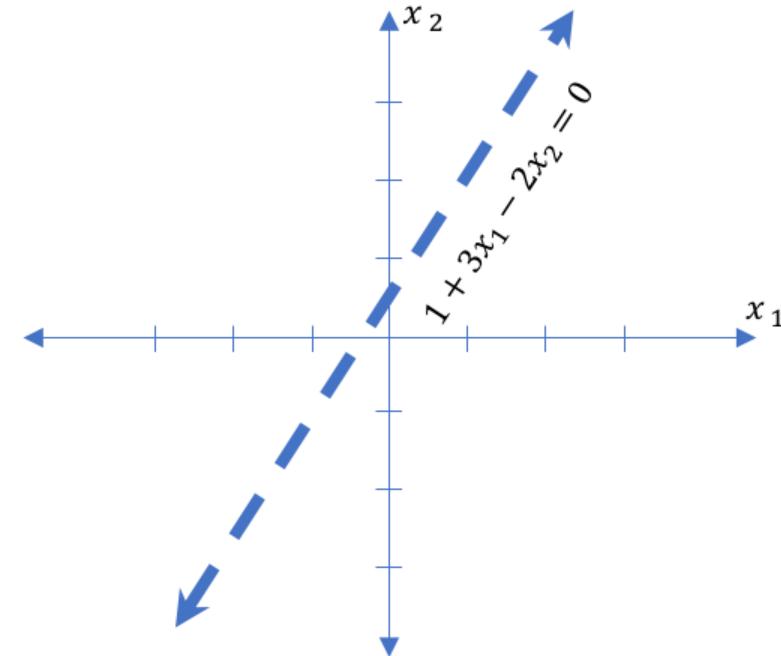
$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{w}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g\left(1 + 3x_1 - 2x_2\right)\end{aligned}$$

This is just a line in 2D!

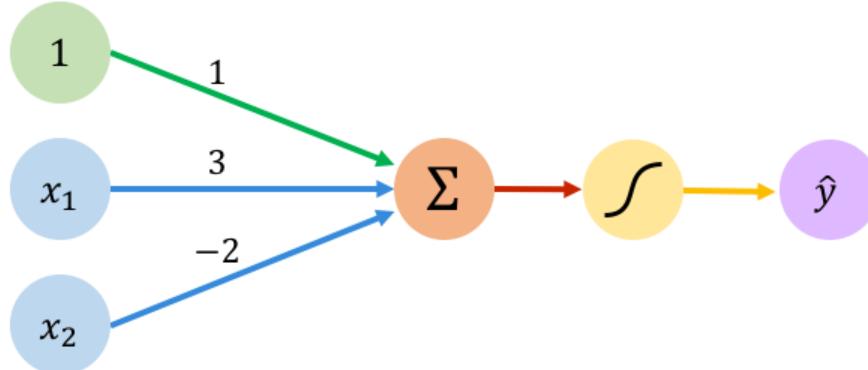
The perceptron : Example



$$\hat{y} = g(1 + 3x_1 - 2x_2)$$



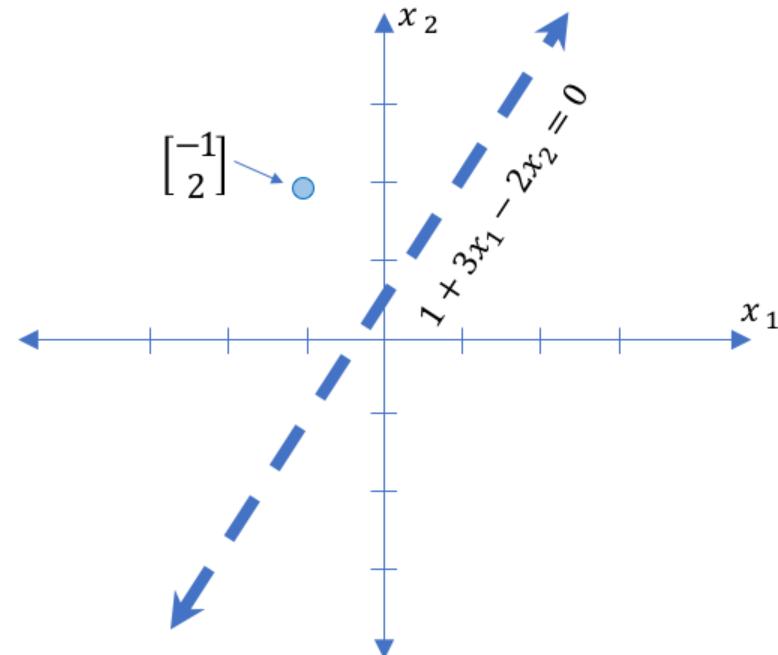
The perceptron : Example



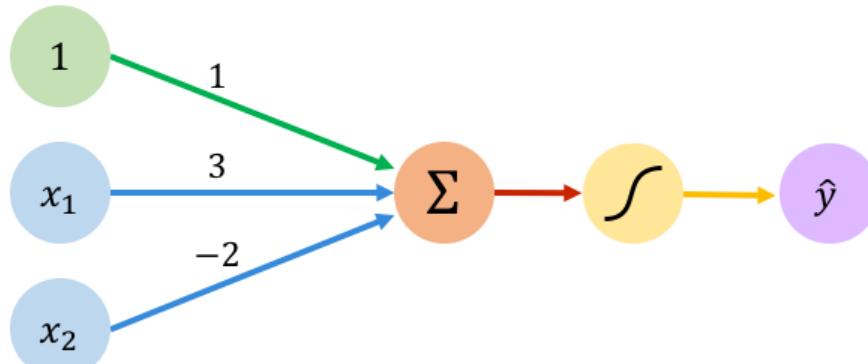
Assume we have input: $\mathbf{X} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$

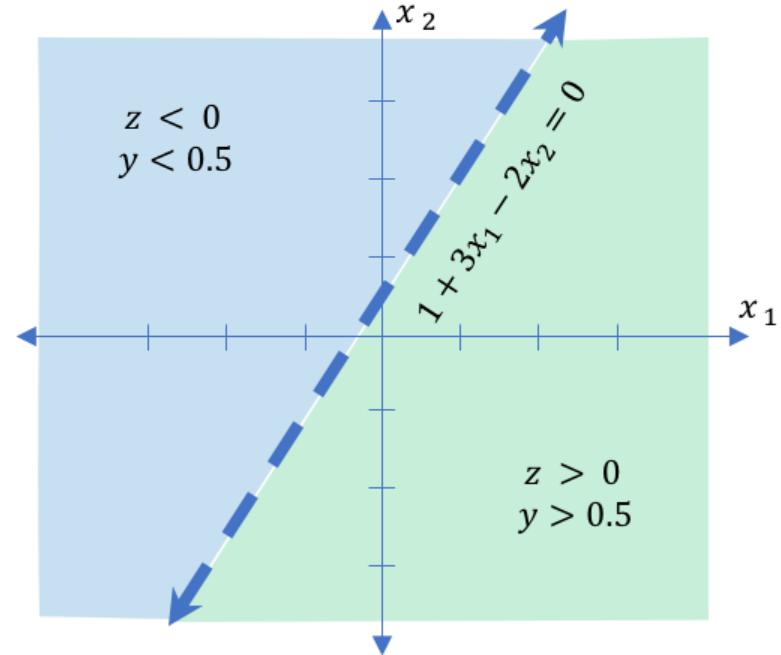
$$\hat{y} = g(1 + 3x_1 - 2x_2)$$



The perceptron : Example



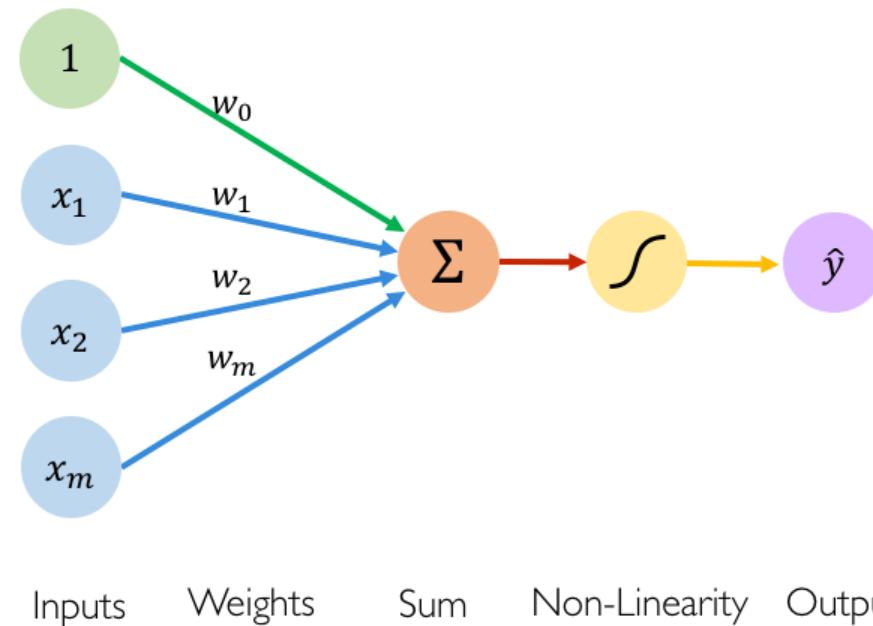
$$\hat{y} = g(1 + 3x_1 - 2x_2)$$



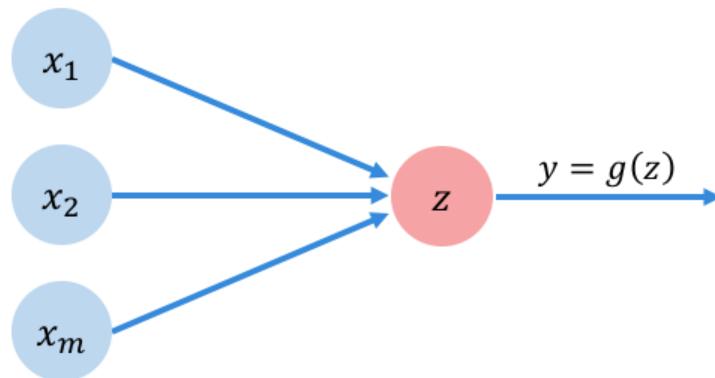
Building Neural Networks with perceptrons

The perceptron : Simplified

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$



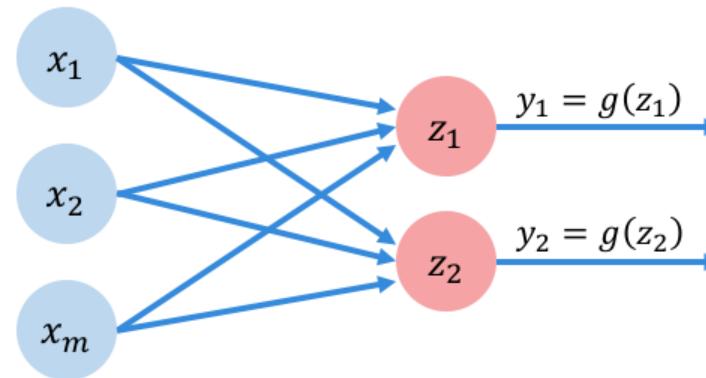
The perceptron : Simplified



$$z = w_0 + \sum_{j=1}^m x_j w_j$$

Multi Output Perceptron

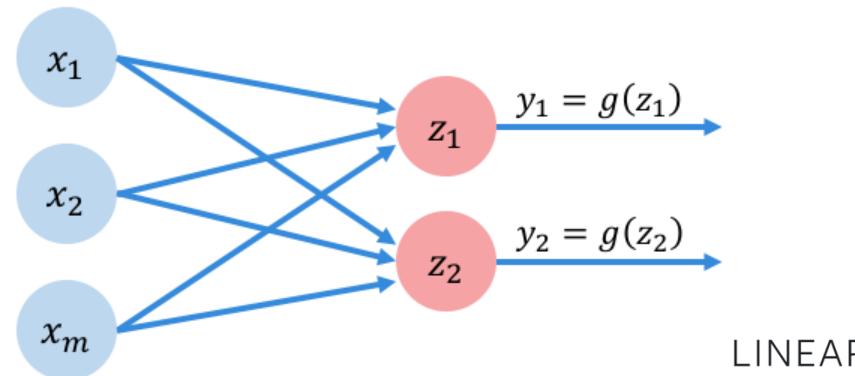
Because all inputs are densely connected to all outputs, these layers are called ***fully connected layers (FC)***



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

Multi Output Perceptron

Because all inputs are densely connected to all outputs, these layers are called ***fully connected layers (FC)***



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

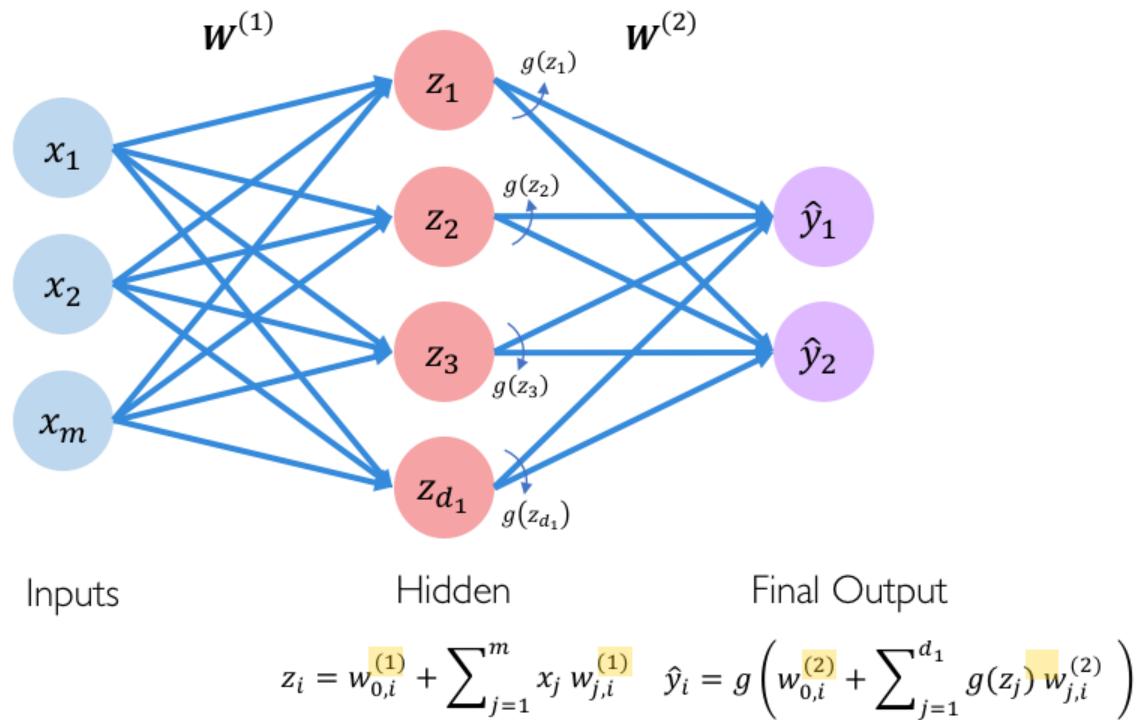
CLASS `torch.nn.Linear(in_features: int, out_features: int, bias: bool = True)`

Applies a linear transformation to the incoming data: $y = xA^T + b$

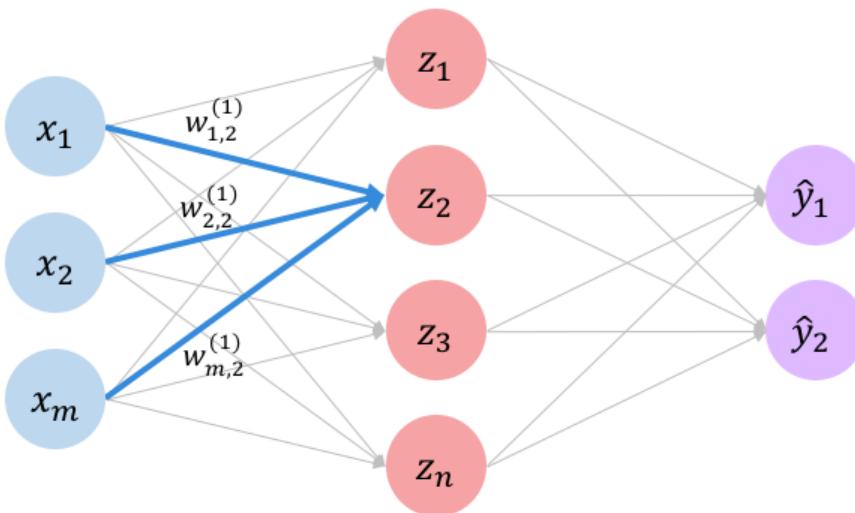
Parameters

- **in_features** – size of each input sample
- **out_features** – size of each output sample
- **bias** – If set to `False`, the layer will not learn an additive bias. Default: `True`

Single Layer Neural Network

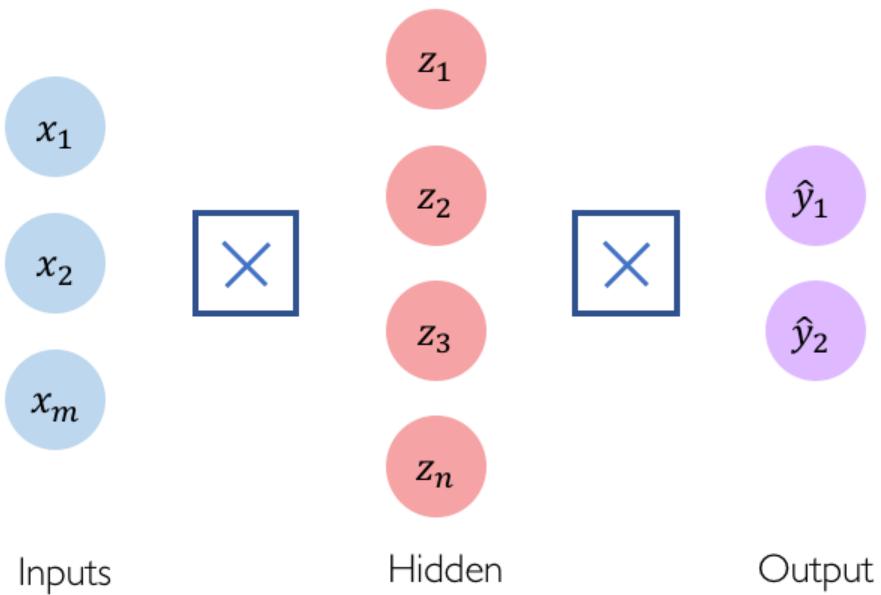


Single Layer Neural Network



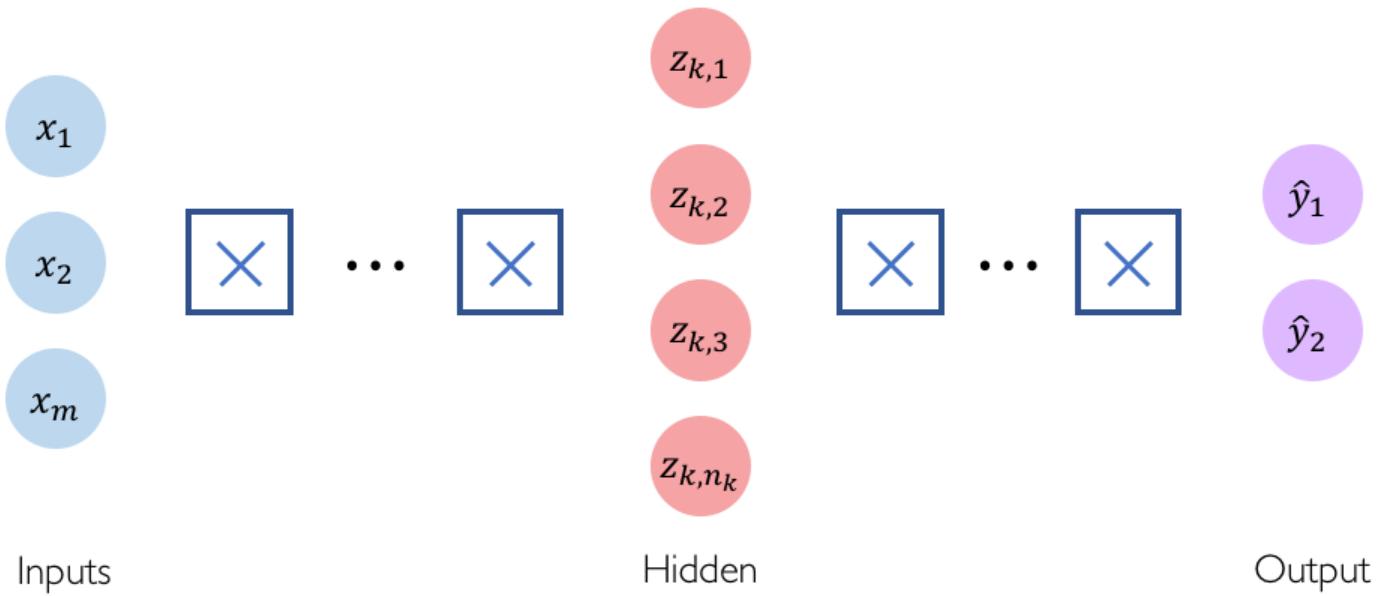
$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$

Single Layer Neural Network



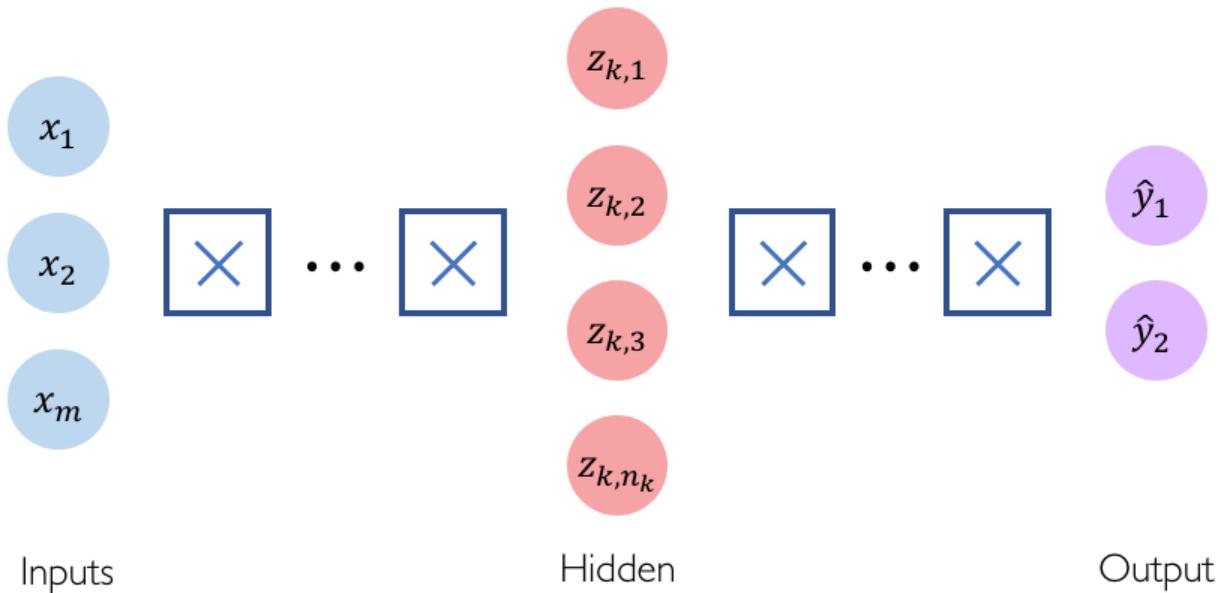
```
import torch  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(m, n),  
    torch.nn.Linear(n, 2)  
)
```

“Deep” Neural Network



$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

“Deep” Neural Network



$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

```
import torch

model = torch.nn.Sequential{
    torch.nn.Linear(m, n1),
    torch.nn.Linear(n1, n2),
    .
    .
    .
    torch.nn.Linear(nk, 2)
}
```

Applying Neural Networks

Example problem

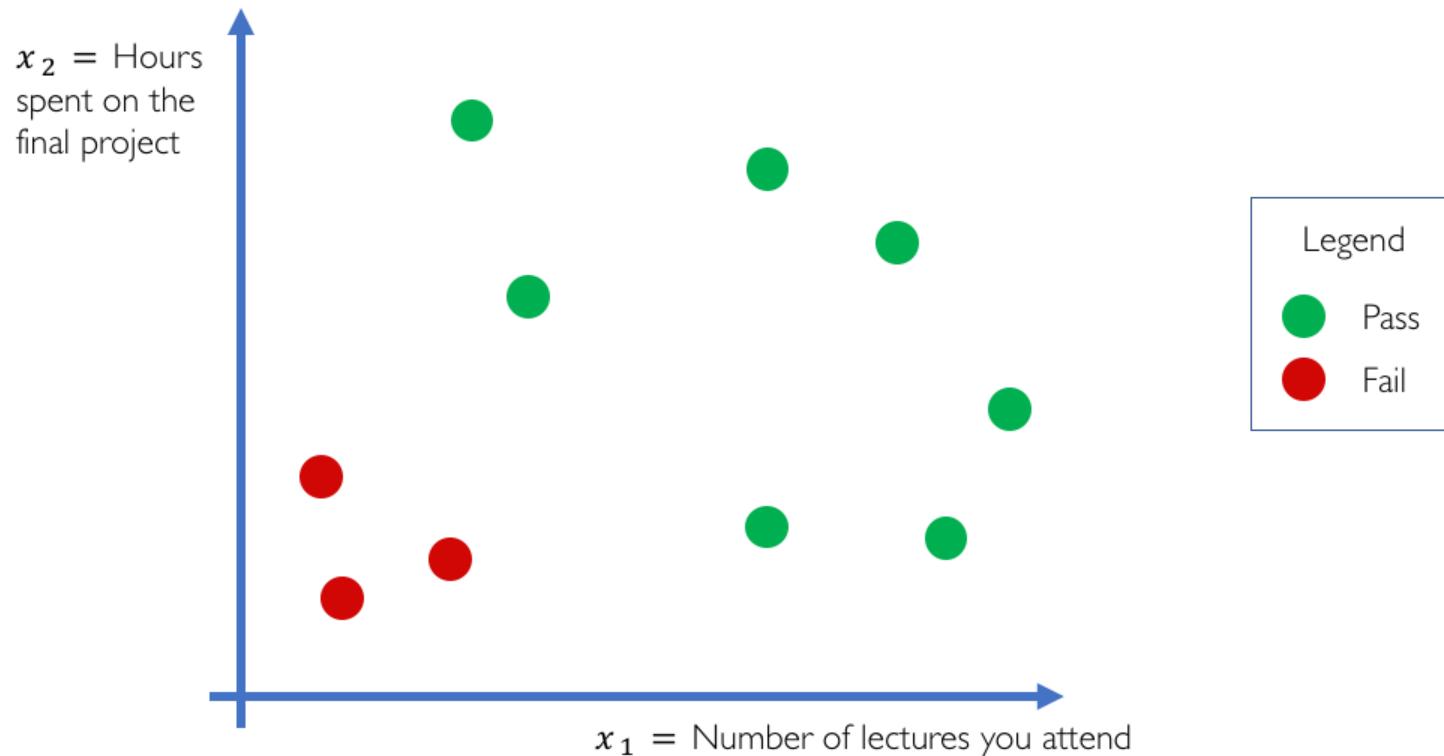
Will I pass this class?

Let's start with a simple two feature model

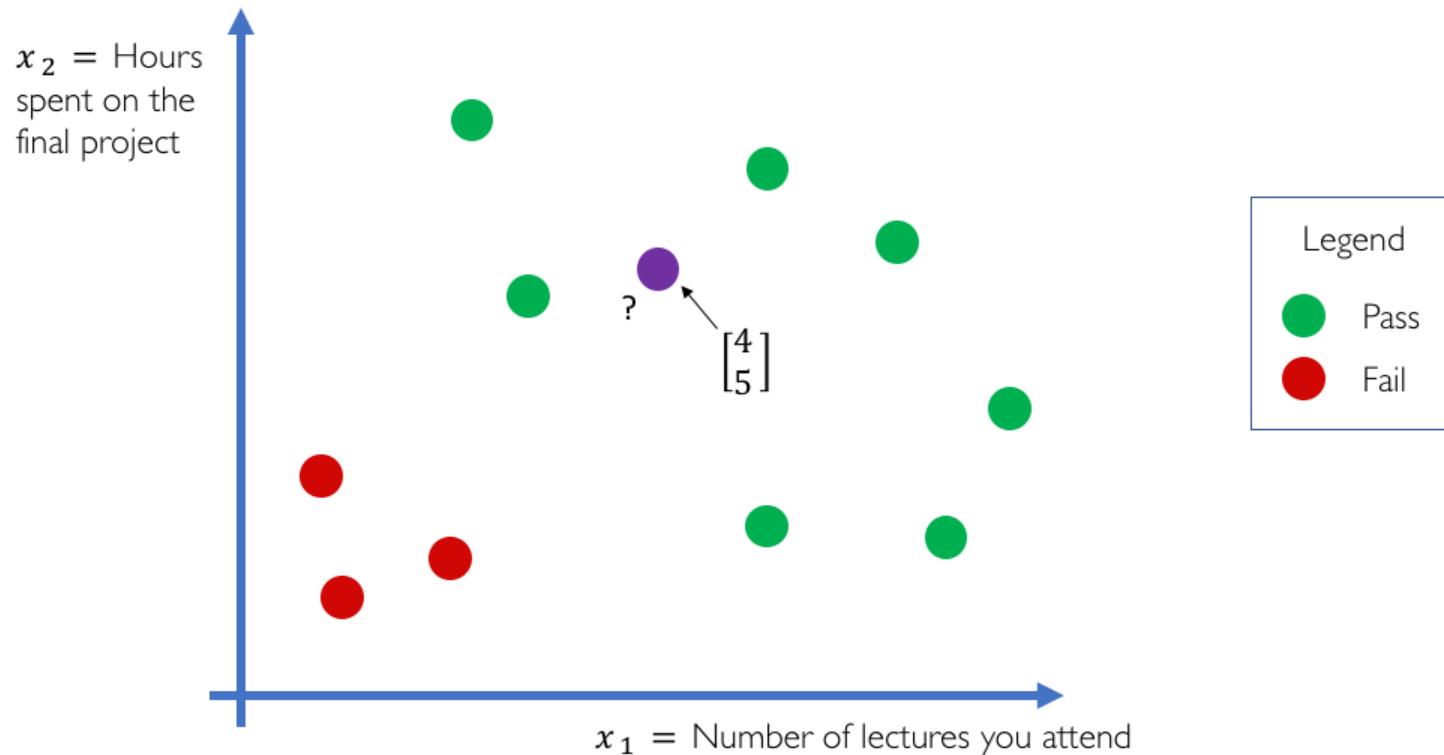
x_1 = Number of lectures you attend

x_2 = Hours spent on the final project

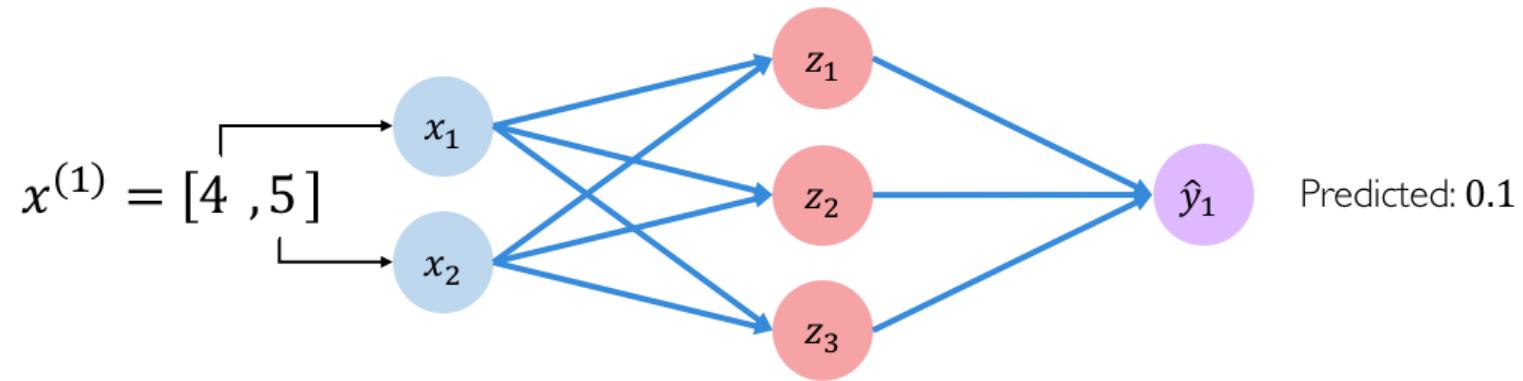
Example problem : Will I pass this class?



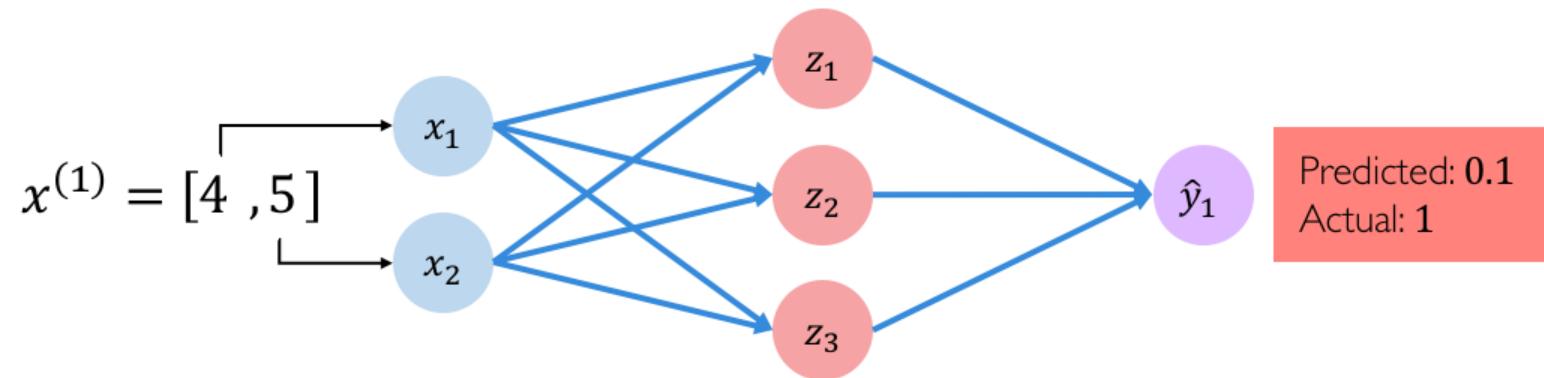
Example problem : Will I pass this class?



Example problem : Will I pass this class?

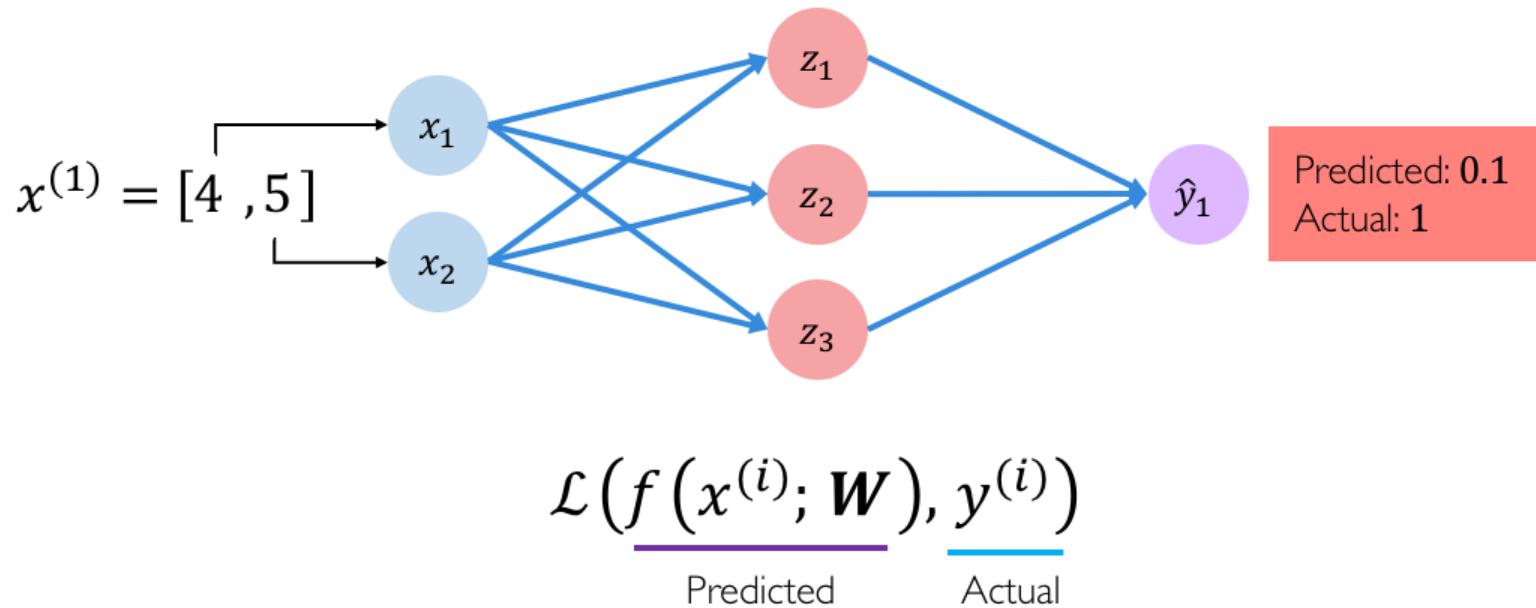


Example problem : Will I pass this class?



Quantifying Loss

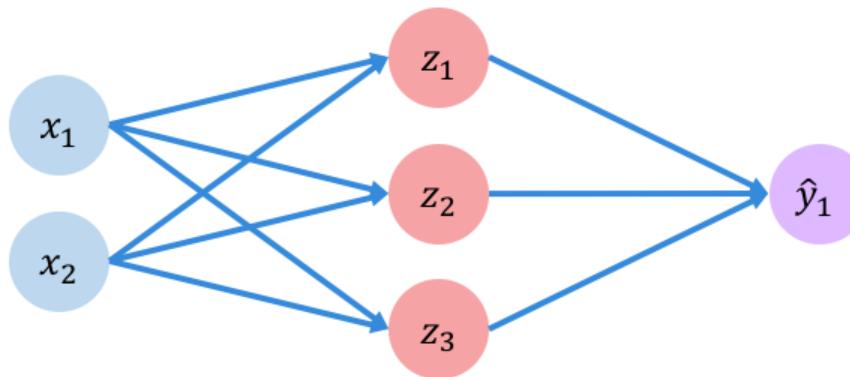
The **loss** of our network measures the cost incurred from incorrect predictions



Quantifying Loss

The **empirical loss** measures the total loss over our entire dataset

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$\begin{array}{c|c} f(x) & y \\ \hline 0.1 & \times | 1 \\ 0.8 & \times | 0 \\ 0.6 & \checkmark | 1 \\ \vdots & \vdots \end{array}$$

- Also known as:
- Objective function
 - Cost function
 - Empirical Risk

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted Actual

Loss Functions



Regression

What is the temperature going to be tomorrow?

PREDICTION
84°

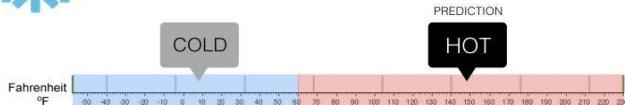


Classification

Will it be Cold or Hot tomorrow?

COLD

PREDICTION
HOT



- Loss function *quantifies* gap between prediction and ground truth
- For regression:
 - Mean Squared Error (MSE)
- For classification:
 - Cross Entropy Loss

Mean Squared Error

$$MSE = \frac{1}{N} \sum (t_i - s_i)^2$$

Prediction
↓
 t_i
Ground Truth
↑
 s_i

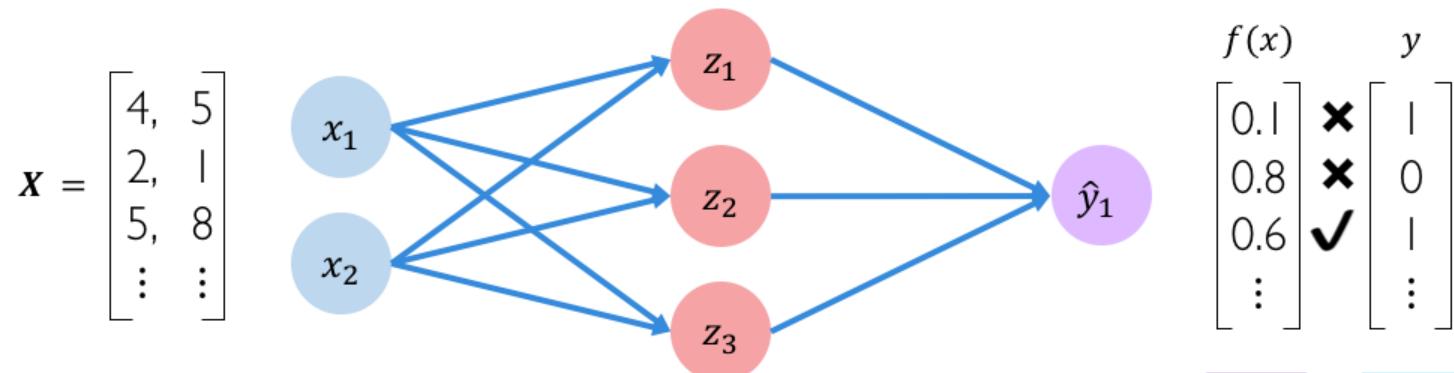
Cross Entropy Loss

$$CE = - \sum_i t_i \log(s_i)$$

Classes
↓
 t_i
Prediction
↓
 s_i
Ground Truth {0,1}

Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \mathbf{W}))}_{\text{Actual}}$$

Predicted Predicted

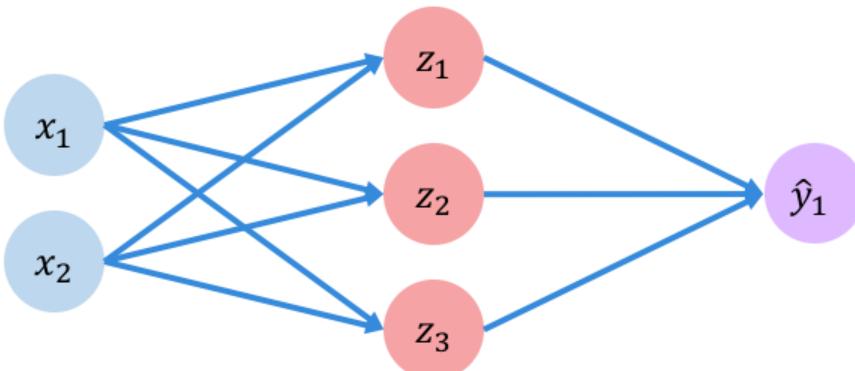
```
loss = F.binary_cross_entropy(predicted, y)
```

```
loss = nn.BCELoss()  
output = loss(predicted, y)
```

Mean Square Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots, & \vdots \end{bmatrix}$$



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \left(\underline{y^{(i)}} - \underline{f(x^{(i)}; \mathbf{W})} \right)^2$$

Actual Predicted

$f(x)$	y
30	✗
80	✗
85	✓
\vdots	\vdots

Final Grades
(percentage)

```
loss = F.mse_loss(predicted, y)
```

```
loss = nn.MSELoss()  
output = loss(predicted, y)
```

Training Neural Networks

Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

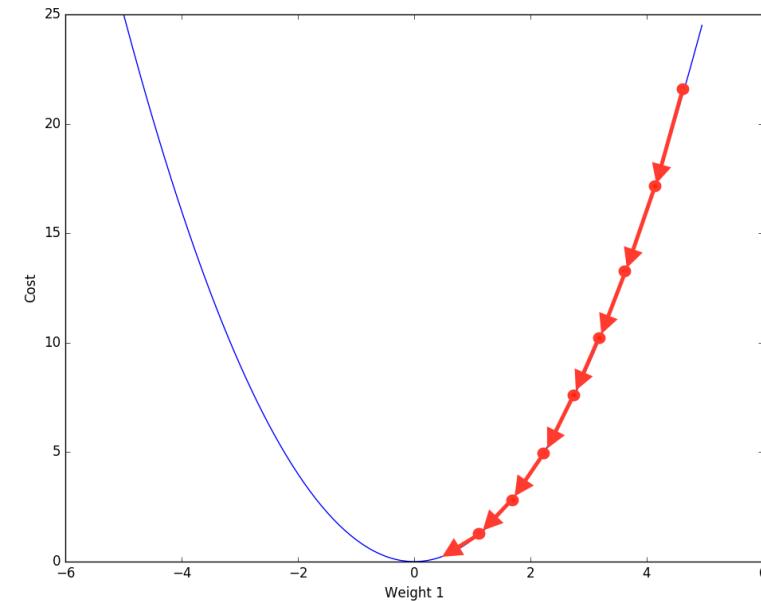
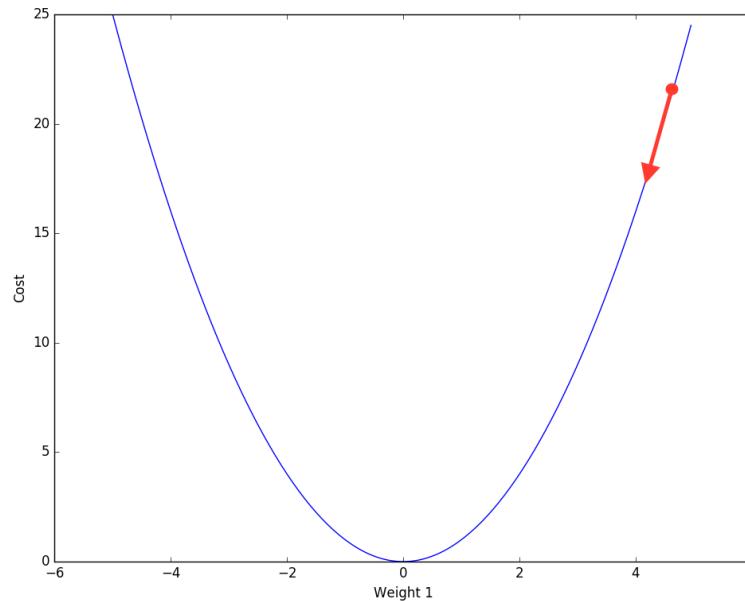
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:
 $\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$

Learning is an Optimization Problem

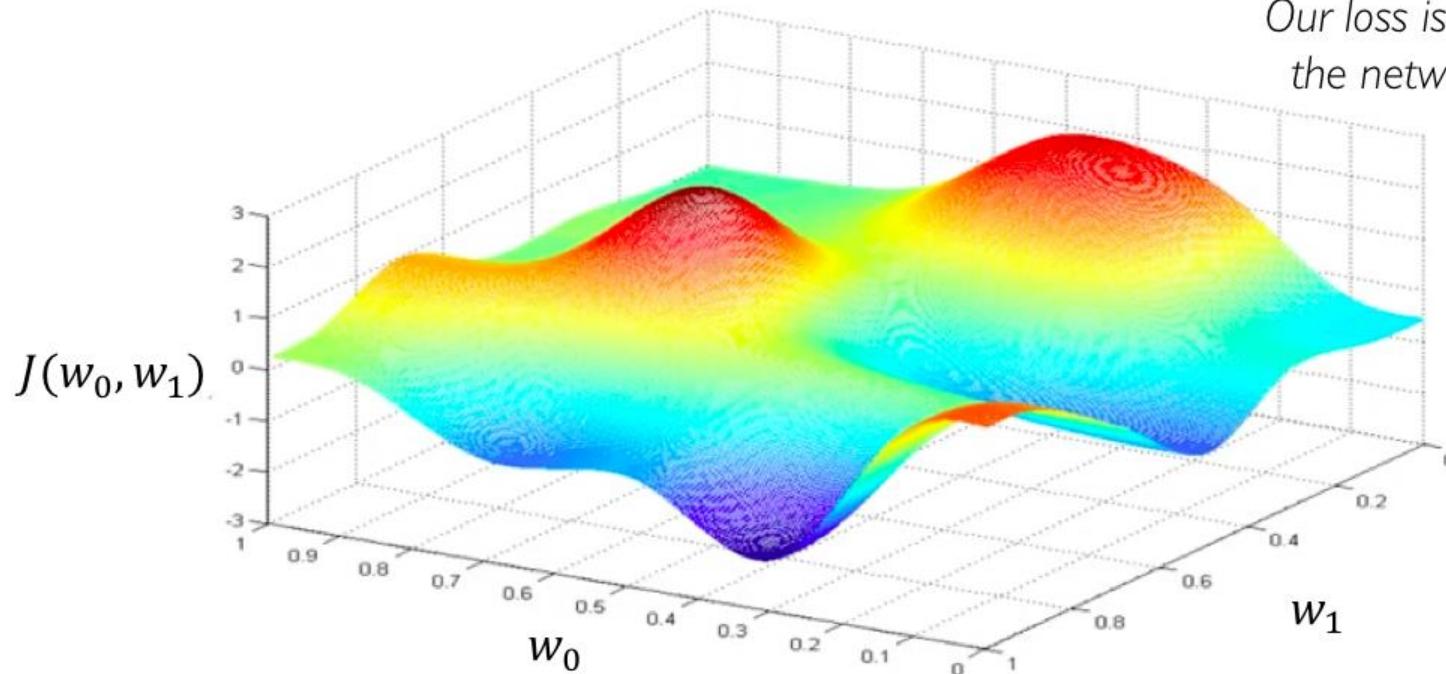
Update the **weights** and **biases** to decrease loss function



SGD: Stochastic Gradient Descent

Loss Optimization

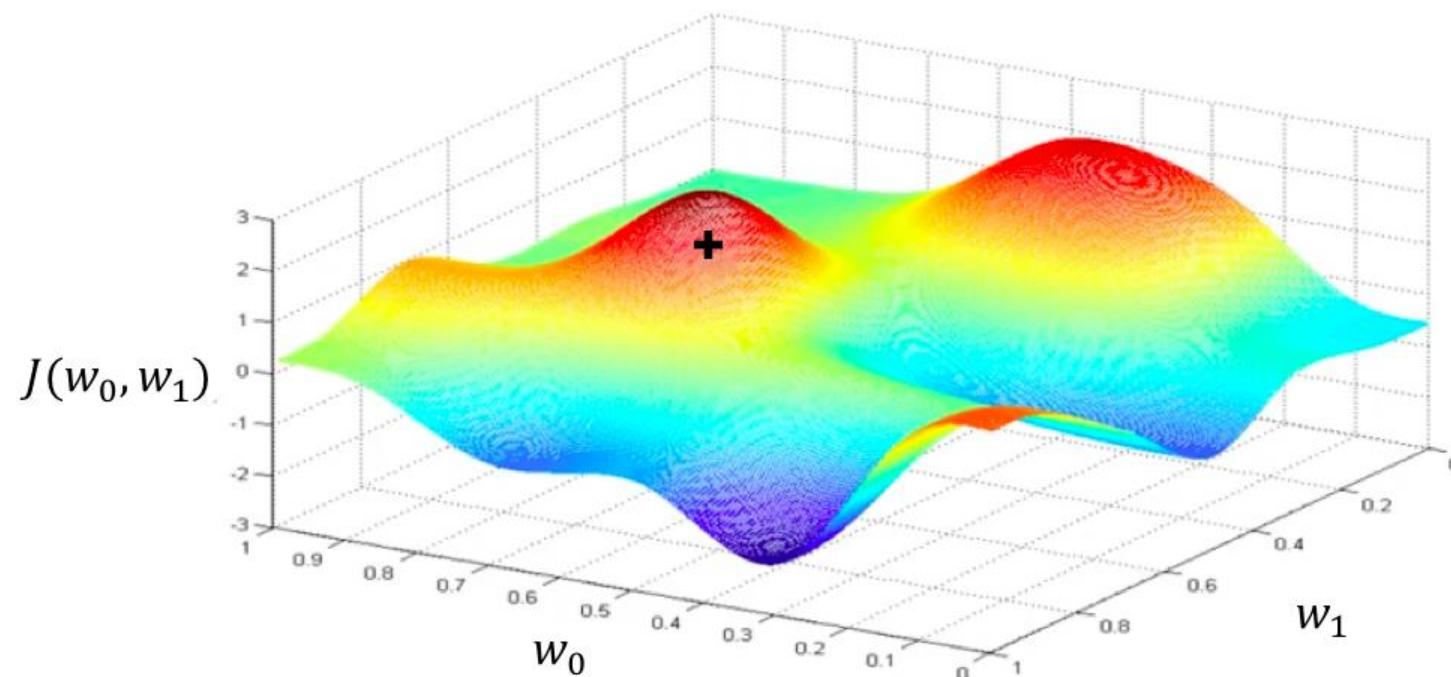
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:
Our loss is a function of
the network weights!

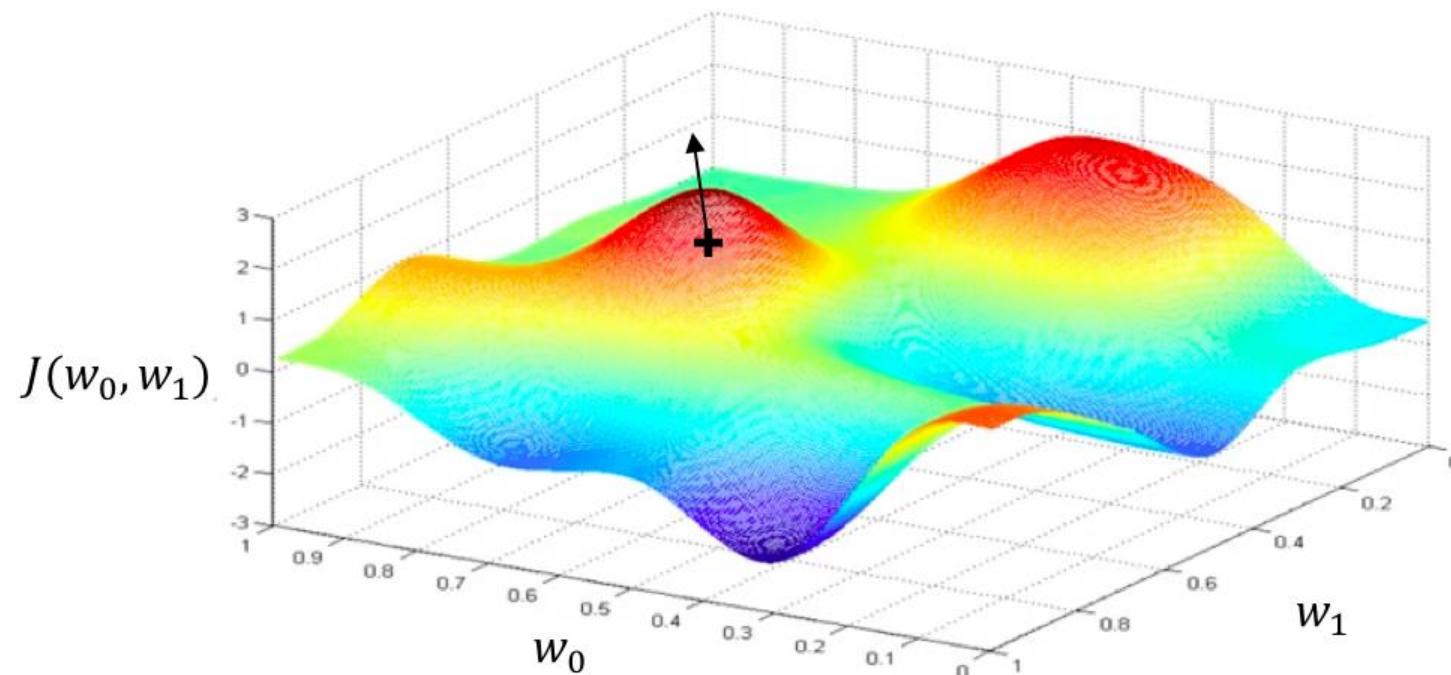
Loss Optimization

Randomly pick an initial (w_0, w_1)



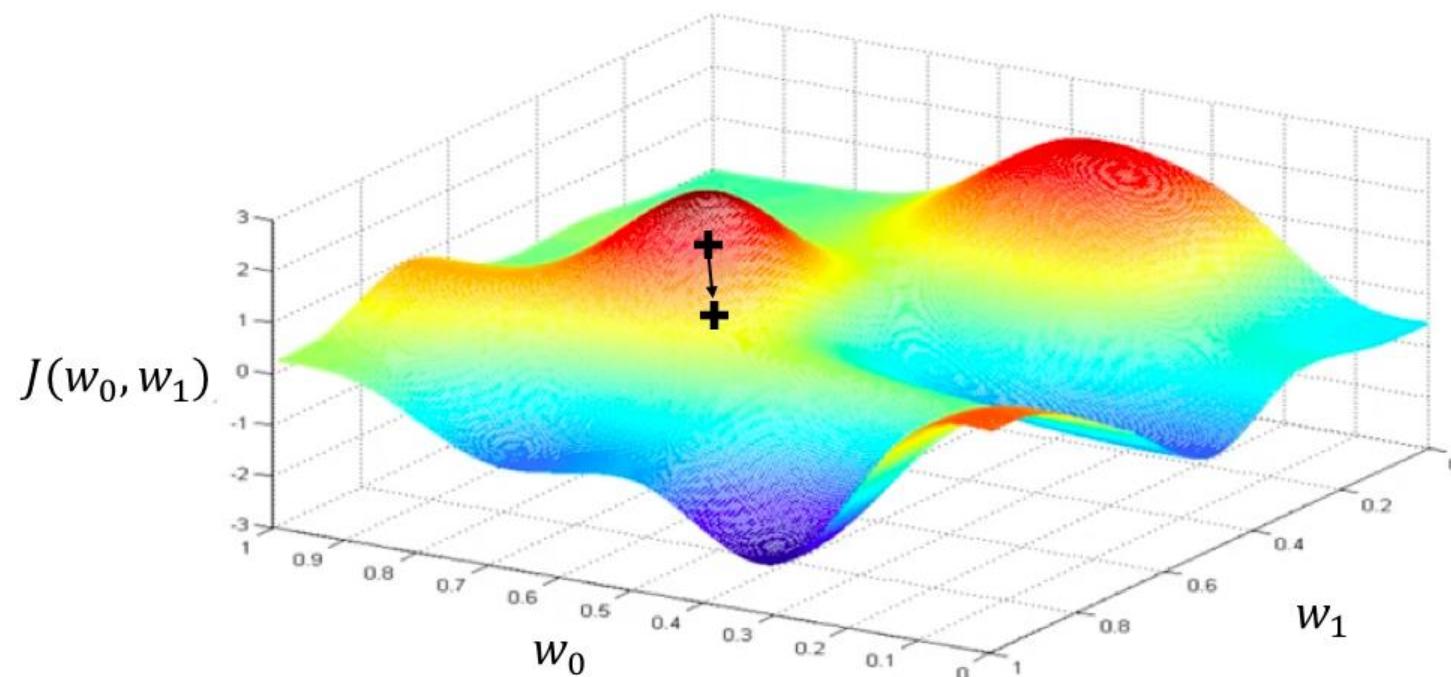
Loss Optimization

Compute gradient, $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$



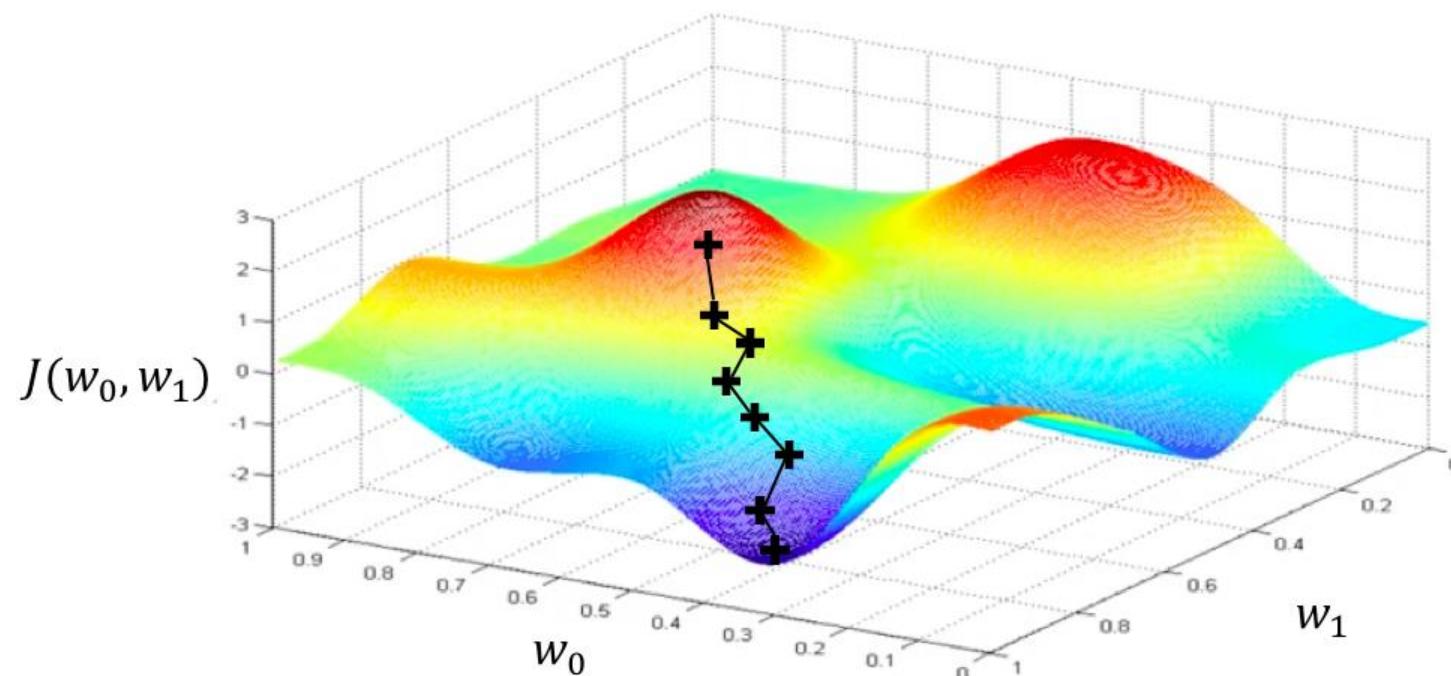
Loss Optimization

Take small step in opposite direction of gradient



Gradient Descent

Repeat until convergence



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

```
import torch

for it in range(max_iter):
    predicted = model(x)
    loss = loss_fn(predicted, y)

    # Zero the gradients before running the backward pass.
    model.zero_grad()

    # Backward pass: compute gradients of the loss
    # with respect to all the learnable parameters of the model.
    loss.backward()

    # Update the weights using gradient descent.
    with torch.no_grad():
        for param in model.parameters():
            param -= learning_rate * param.grad
```

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

```
import torch

for it in range(max_iter):
    predicted = model(x)
    loss = loss_fn(predicted, y)

    # Zero the gradients before running the backward pass.
    model.zero_grad()

    # Backward pass: compute gradients of the loss
    # with respect to all the learnable parameters of the model.
    loss.backward()

    # Update the weights using gradient descent.
    with torch.no_grad():
        for param in model.parameters():
            param -= learning_rate * param.grad
```

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

```
import torch

for it in range(max_iter):
    predicted = model(x)
    loss = loss_fn(predicted, y)

    # Zero the gradients before running the backward pass.
    model.zero_grad()

    # Backward pass: compute gradients of the loss
    # with respect to all the learnable parameters of the model.
    loss.backward()

    # Update the weights using gradient descent.
    with torch.no_grad():
        for param in model.parameters():
            param -= learning_rate * param.grad
```

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

```
import torch

optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)
for it in range(max_iter):

    predicted = model(x)

    loss = loss_fn(predicted, y)

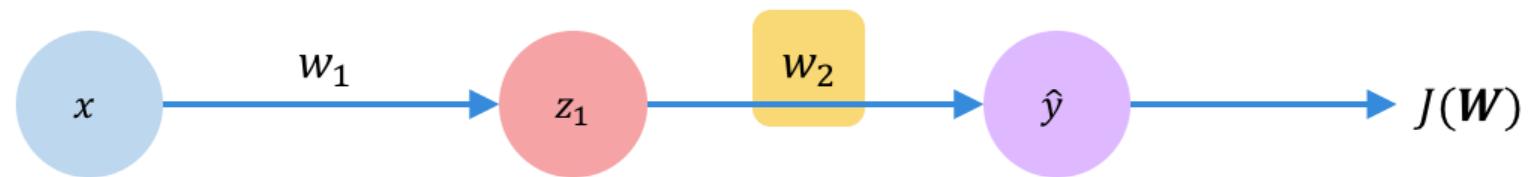
    # Zero the gradients before running the backward pass.
    # model.zero_grad()
    optimizer.zero_grad()

    # Backward pass: compute gradients of the loss
    # with respect to all the learnable parameters of the model.
    loss.backward()

    # Update the weights using gradient descent.
    with torch.no_grad():
        # for param in model.parameters():
        #     param -= learning_rate * param.grad
    optimizer.step()
```

Computing Gradient : Backpropagation

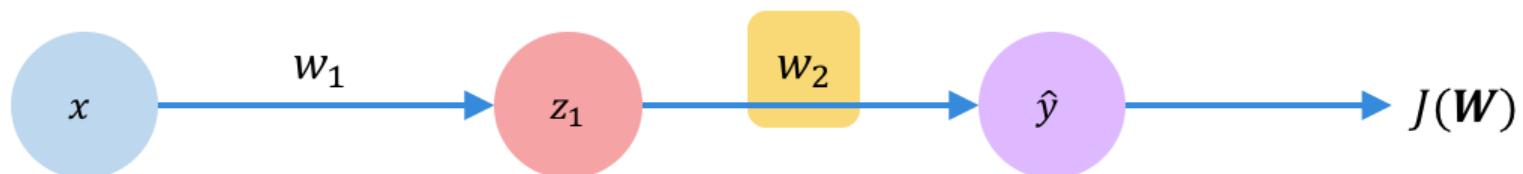
(aka. backprop.)



How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

Computing Gradient : Backpropagation

(aka. backprop.)

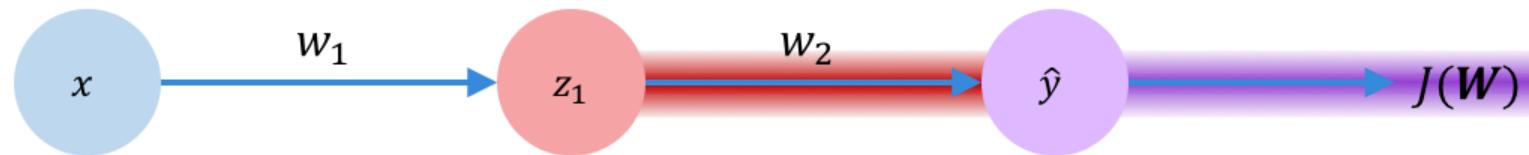


$$\frac{\partial J(\mathbf{W})}{\partial w_2} =$$

Let's use the chain rule!

Computing Gradient : Backpropagation

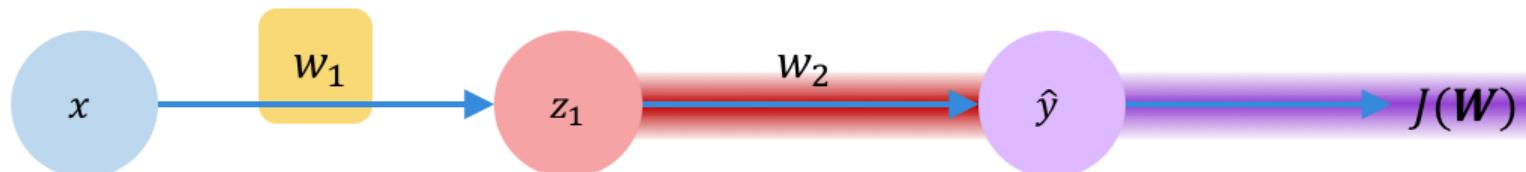
(aka. backprop.)



$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} * \overline{\frac{\partial \hat{y}}{\partial w_2}}$$

Computing Gradient : Backpropagation

(aka. backprop.)



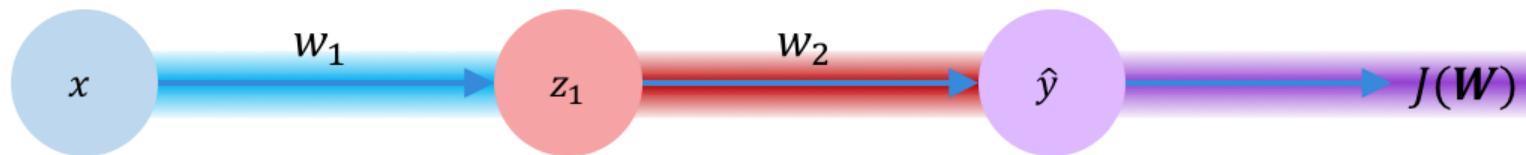
$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule!

Apply chain rule!

Computing Gradient : Backpropagation

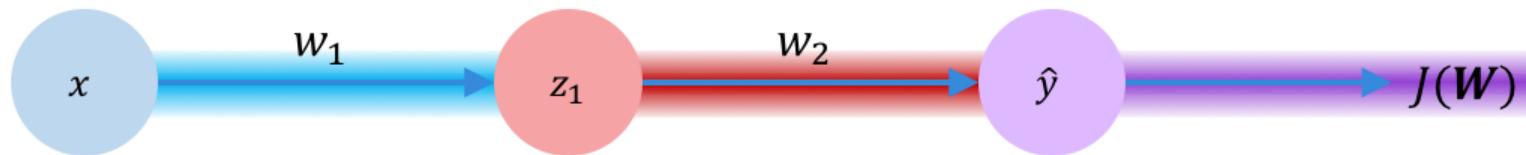
(aka. backprop.)



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} * \underline{\frac{\partial \hat{y}}{\partial z_1}} * \underline{\frac{\partial z_1}{\partial w_1}}$$

Computing Gradient : Backpropagation

(aka. backprop.)



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

Repeat this for **every weight in the network** using gradients from later layers

Back propagation in practice

Computer does backprop. for you

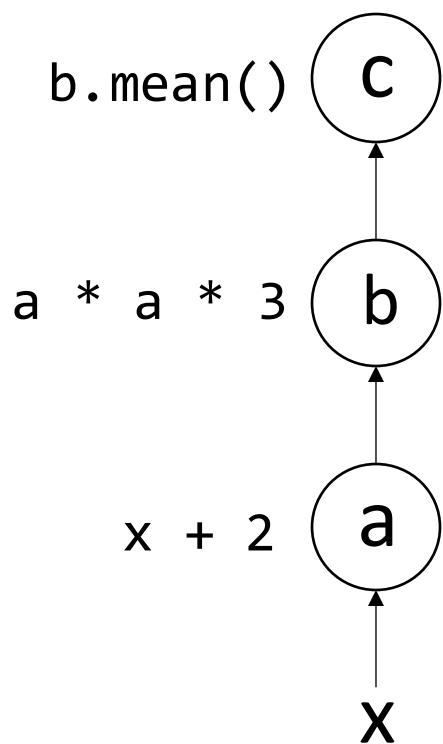
- Autograd
- framework: Tensorflow, Pytorch, MXNet

A neural network model is represented as Directed Acyclic Graph (DAG)

- Each node contains a mathematical operation.
- Each node contains the derivative of the math operation.
- The error signal is backpropagated from the output nodes **to the input** nodes

Math program as DAG

```
x : [[1,1], [1,1]]  
a = x + 2  
b = a * a * 3  
c = b.mean()
```

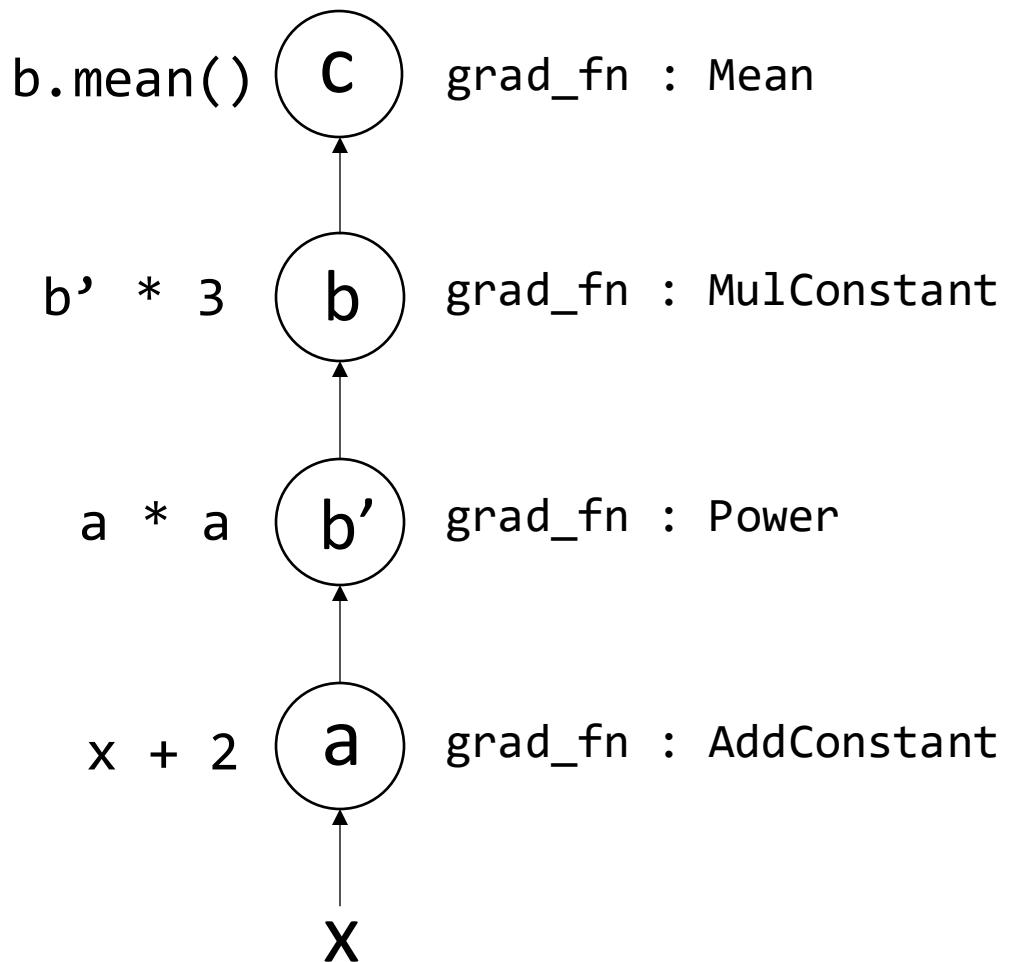


Math program as DAG

```
x : [[1,1], [1,1]]  
a = x + 2  
b' = a * a  
b = b' * 3  
c = b.mean()
```

Each **grad_fn** has

- **forward**
 - Compute forward computation
 - Save input & misc. for backward
- **backward**
 - given δ^l , calculate δ^{l-1}

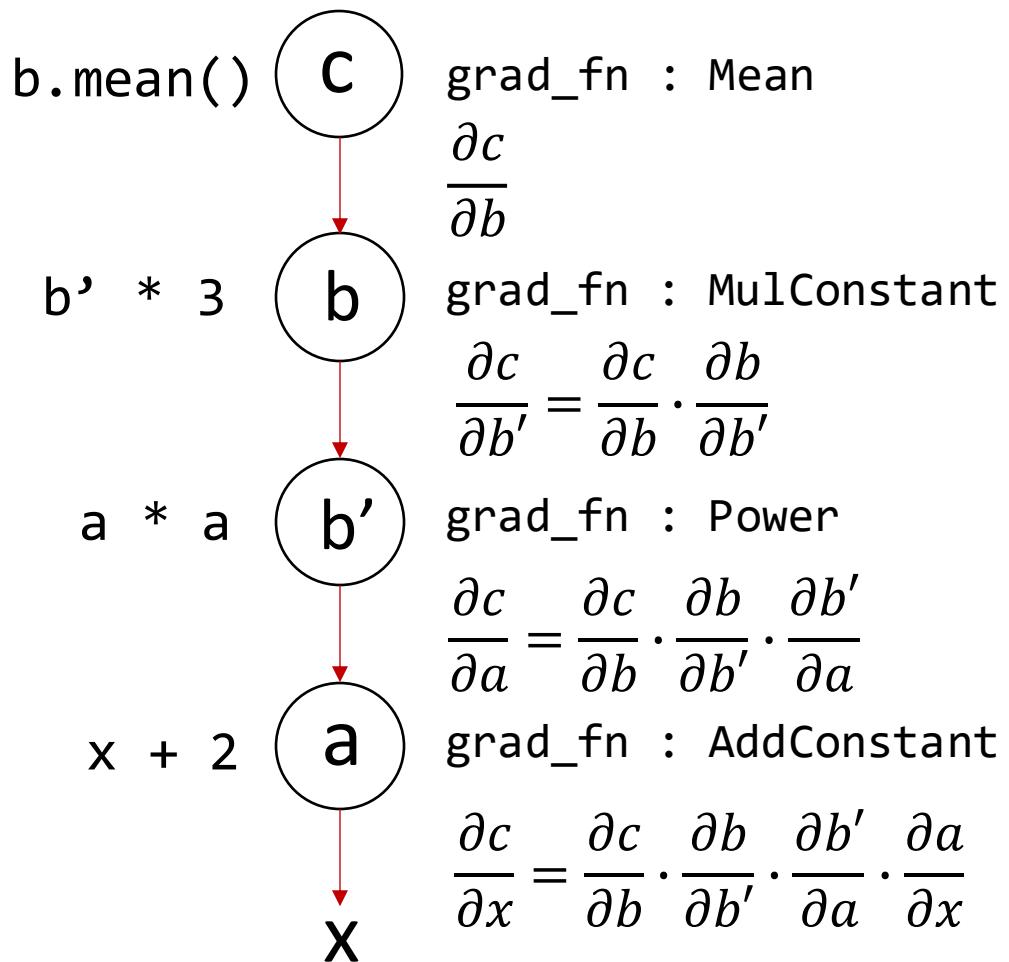


Math program as DAG

```
x : [[1,1], [1,1]]  
a = x + 2  
b' = a * a  
b = b' * 3  
c = b.mean()
```

Each **grad_fn** has

- **forward**
 - Compute forward computation
 - Save input & misc. for backward
- **backward**
 - given δ^l , calculate δ^{l-1}

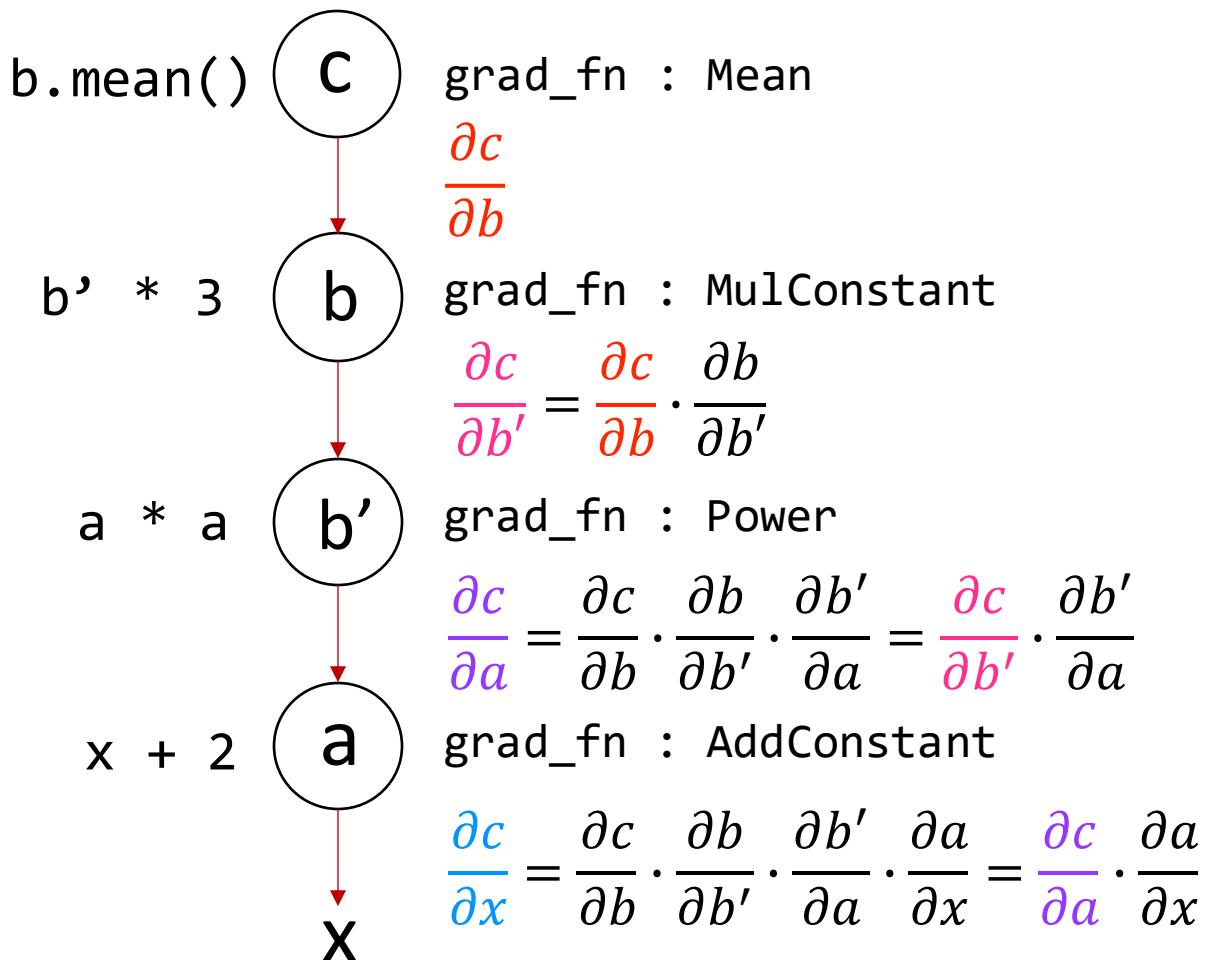


Math program as DAG

```
x : [[1,1], [1,1]]  
a = x + 2  
b' = a * a  
b = b' * 3  
c = b.mean()
```

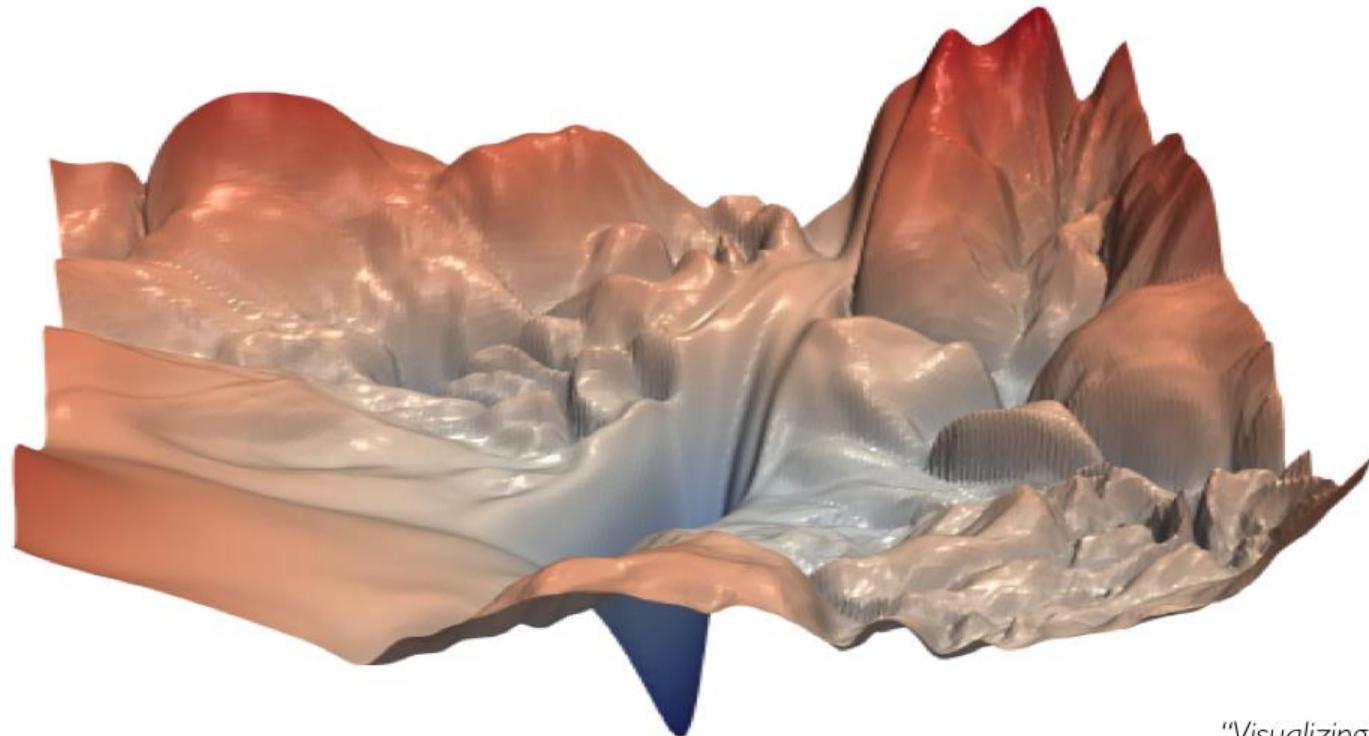
Each **grad_fn** has

- **forward**
 - Compute forward computation
 - Save input & misc. for backward
- **backward**
 - given δ^l , calculate δ^{l-1}



Neural Networks in Practice: Optimization

Training Neural Network is difficult



*"Visualizing the loss landscape
of neural nets". Dec 2017.*

Loss Functions Can Be Difficult to Optimize

Remember:

Optimization through gradient descent

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

Loss Functions Can Be Difficult to Optimize

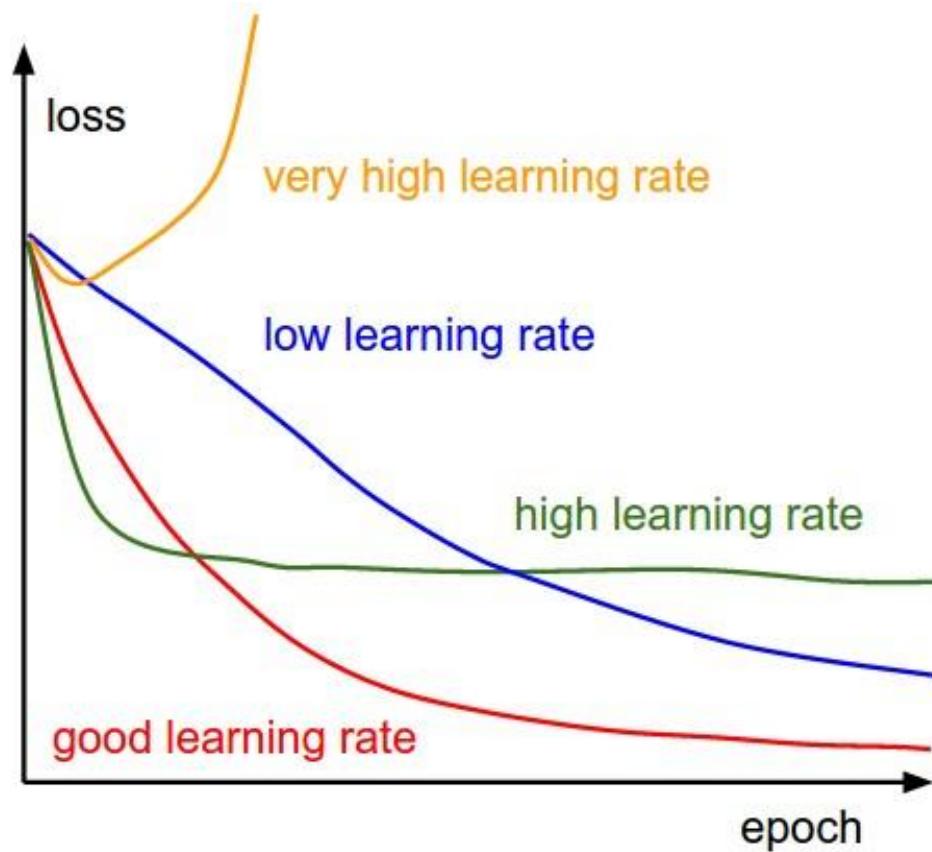
Remember:

Optimization through gradient descent

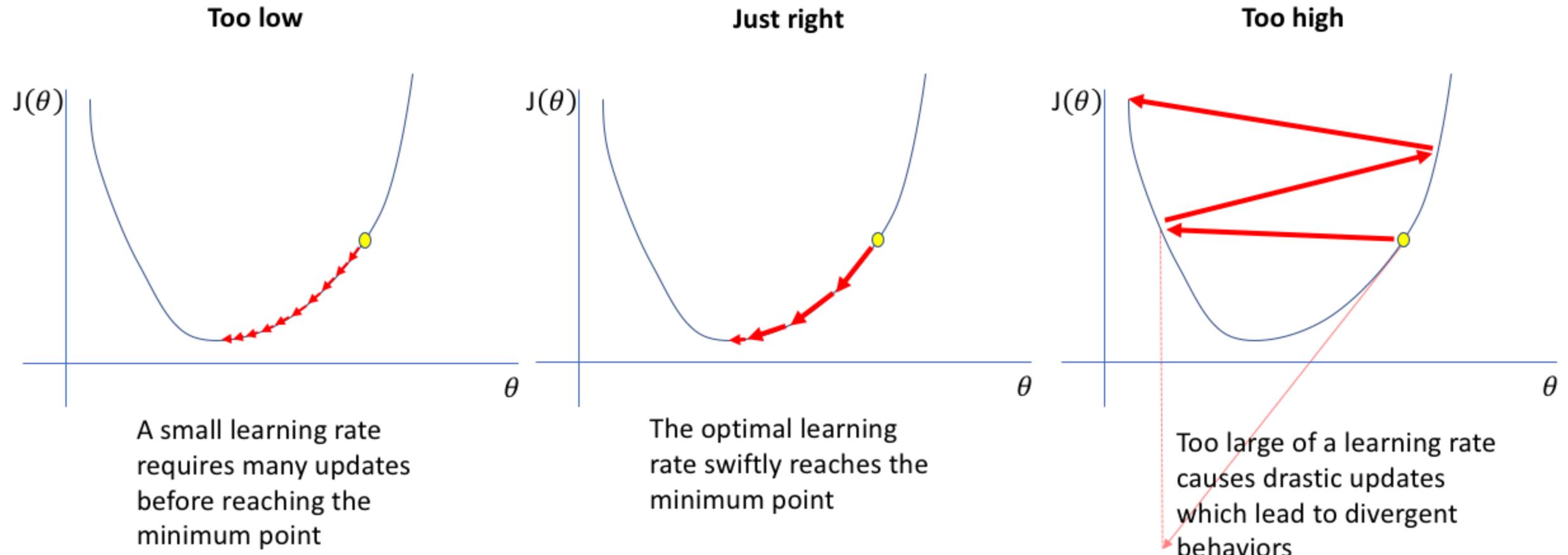
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

How can we set the
learning rate?

Setting the Learning Rate



Setting the Learning Rate



How to deal with this?

Idea I:

Try lots of different learning rates and see what works “just right”

How to deal with this?

Idea 1:

Try lots of different learning rates and see what works “just right”

Idea 2:

Do something smarter!

Design an adaptive learning rate that “adapts” to the landscape

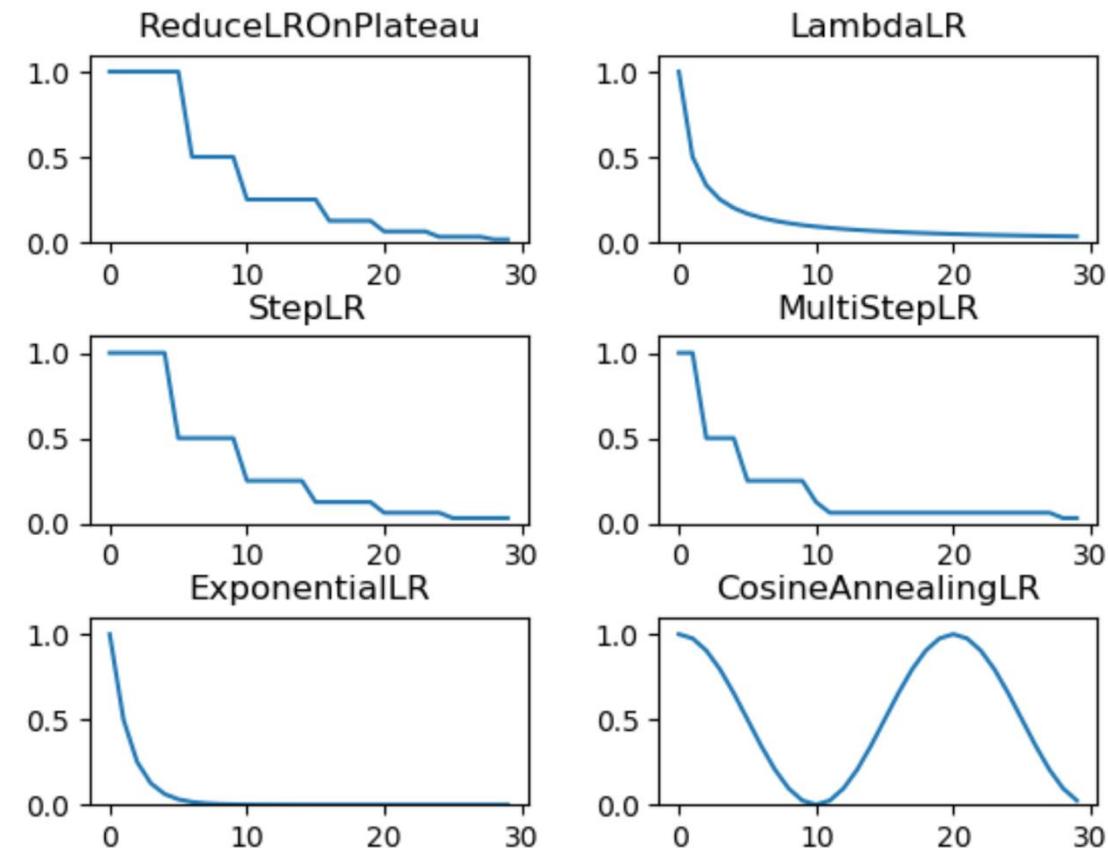
Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - how large gradient is
 - how fast learning is happening
 - size of particular weights
 - etc...

Adaptive Learning Rates

```
import torch
```

```
torch.optim.lr_scheduler.ReduceLROnPlateau(...)  
optim.lr_scheduler.LambdaLR(...)  
optim.lr_scheduler.StepLR(...)  
optim.lr_scheduler.MultiStepLR(...)  
optim.lr_scheduler.ExponentialLR(...)  
optim.lr_scheduler.CosineAnnealingLR(...)
```



Optimizers

Algorithm

- SGD
- Adam
- Adadelta
- Adagrad
- RMSProp



Pytorch Implementation

```
torch.optim.SGD(...)  
torch.optim.Adam(...)  
torch.optim.Adadelta(...)  
torch.optim.Adagrad(...)  
torch.optim.RMSProp(...)
```

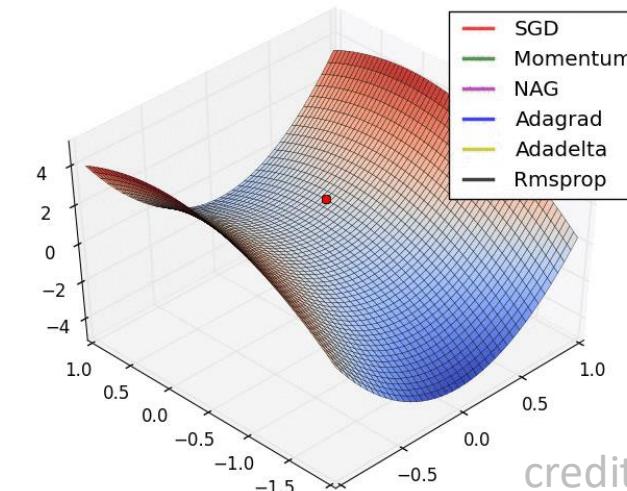
Reference

Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952.

Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.

Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.

Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.



credit : CS231n

Training Process

```
import torch

#define model
model = torch.nn.Sequential(...)

#initial learning rate
learning_rate = 0.05

# optimizer
optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)

#lr scheduler
# lr = 0.05      if epoch < 30
# lr = 0.005     if 30 <= epoch < 80
# lr = 0.0005    if epoch >= 80
scheduler = torch.optim.lr_scheduler.MultiStepLR(optimizer,
                                                milestones=[30,80])

for it in range(max_iter):

    predicted = model(x)

    loss = loss_fn(predicted, y)

    # Zero the gradients before running the backward pass.
    optimizer.zero_grad()

    # Backward pass: compute gradient of the loss
    # with respect to all the learnable parameters of the model.
    loss.backward()

    # Update the weights using gradient descent.
    optimizer.step()

    # adjust learning_rate
    scheduler.step()
```

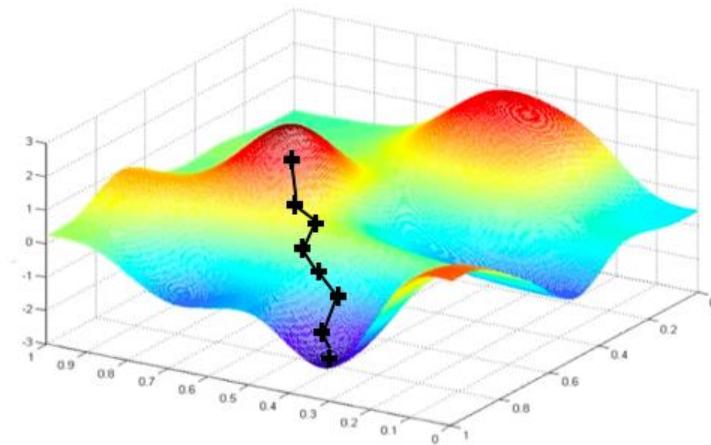
Neural Networks in Practice:

Mini-batch

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

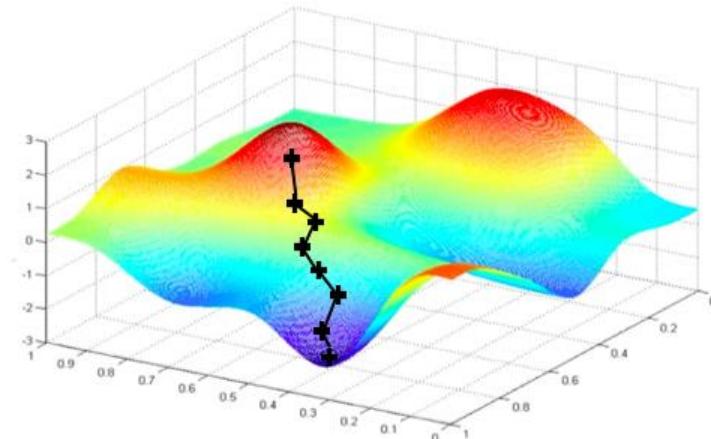


Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

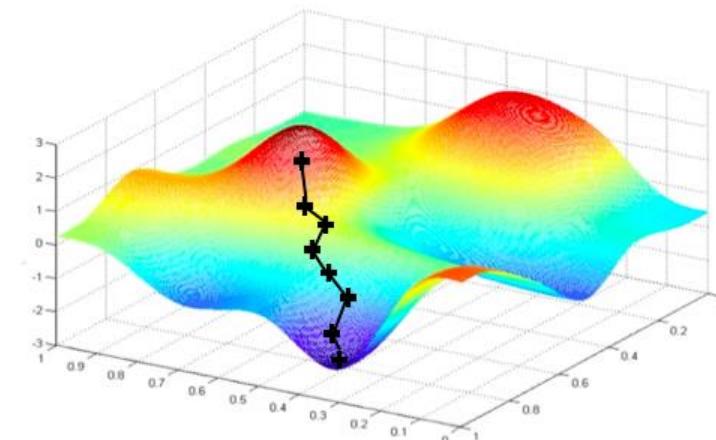
Can be very
computationally
intensive to compute!



Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{w})}{\partial \mathbf{w}}$
5. Update weights, $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$
6. Return weights

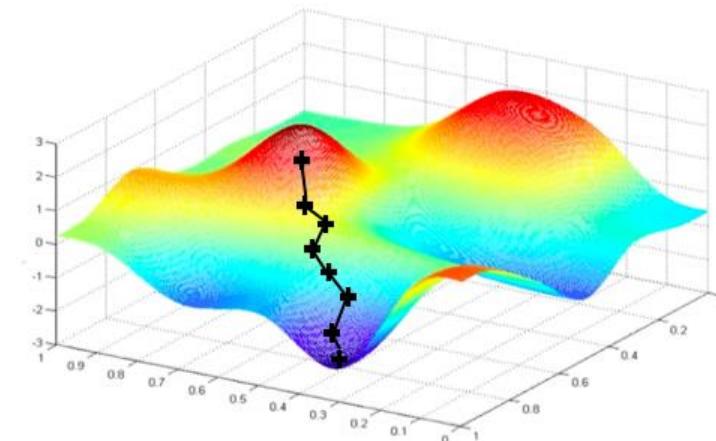


Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

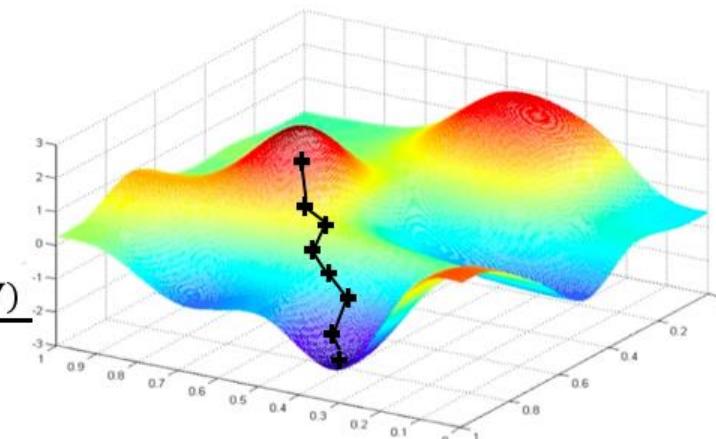
Easy to compute but
very noisy (stochastic)!



Stochastic Gradient Descent

Algorithm

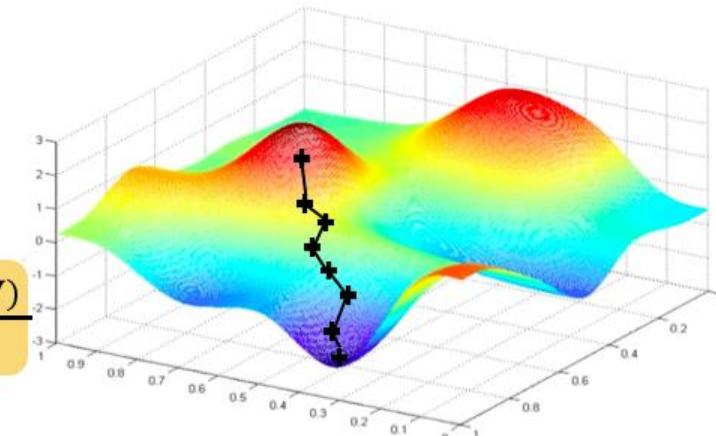
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$
5. Update weights,
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$
6. Return weights



Fast to compute and a much better estimate of the true gradient!

Mini-Batches while Training

More accurate estimation of gradient

Smoother convergence

Allows for larger learning rates

Mini-Batches while Training

More accurate estimation of gradient

Smoother convergence

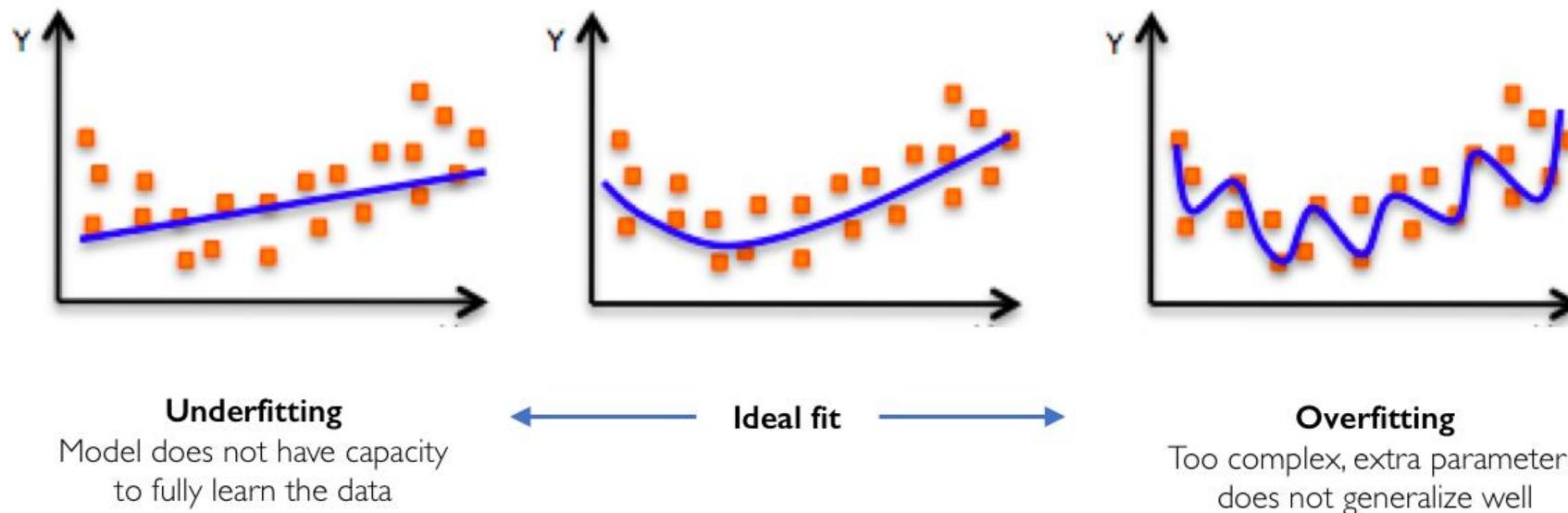
Allows for larger learning rates

Mini-batches lead to fast training!

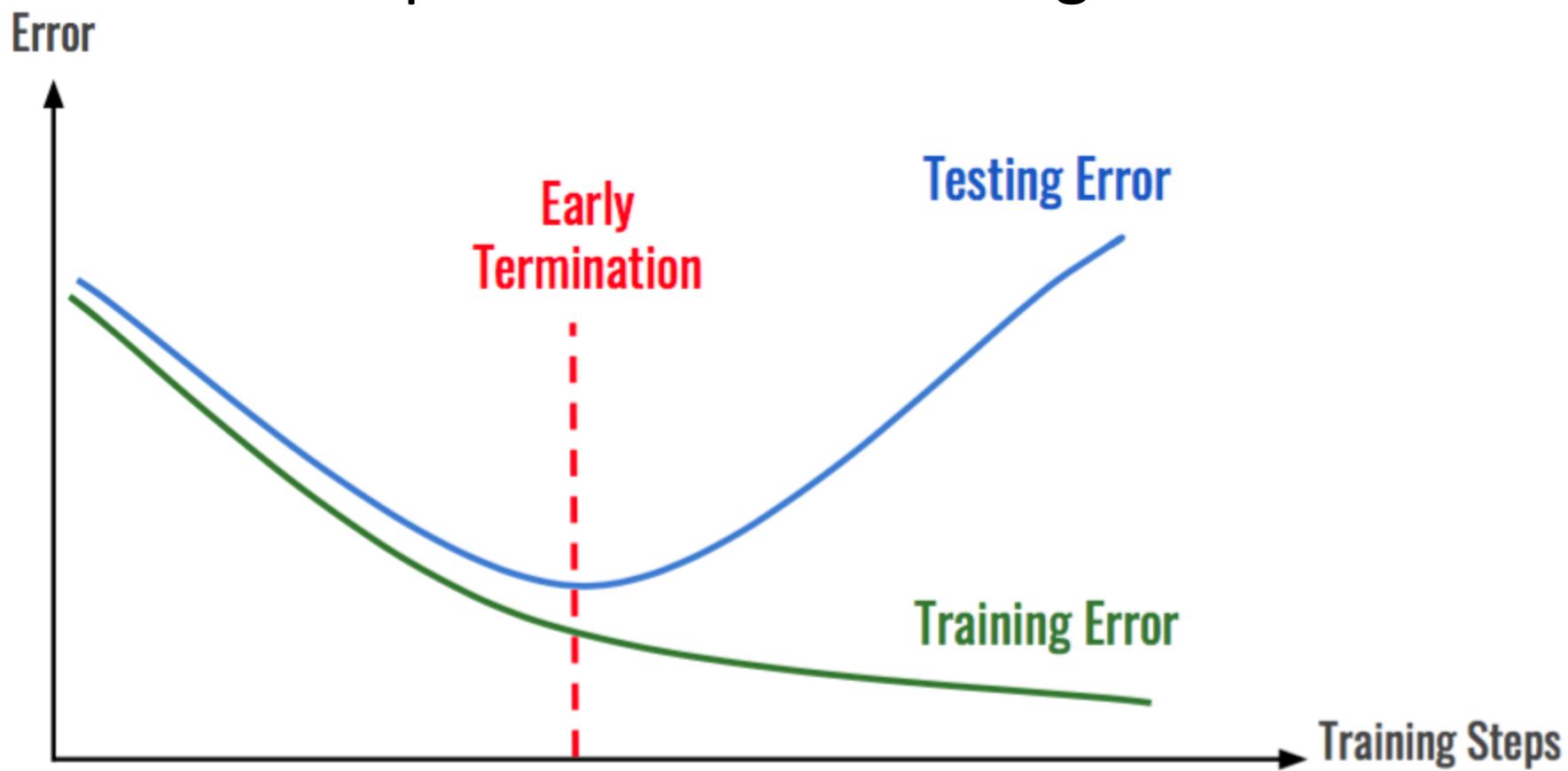
Can parallelize computation + achieve significant speed increases on GPU's

Neural Networks in Practice: Overfitting

The problem of Overfitting



The problem of Overfitting



Regularization

What is it?

Technique that **constrains** our optimization problem to **discourage** complex models

Regularization

What is it?

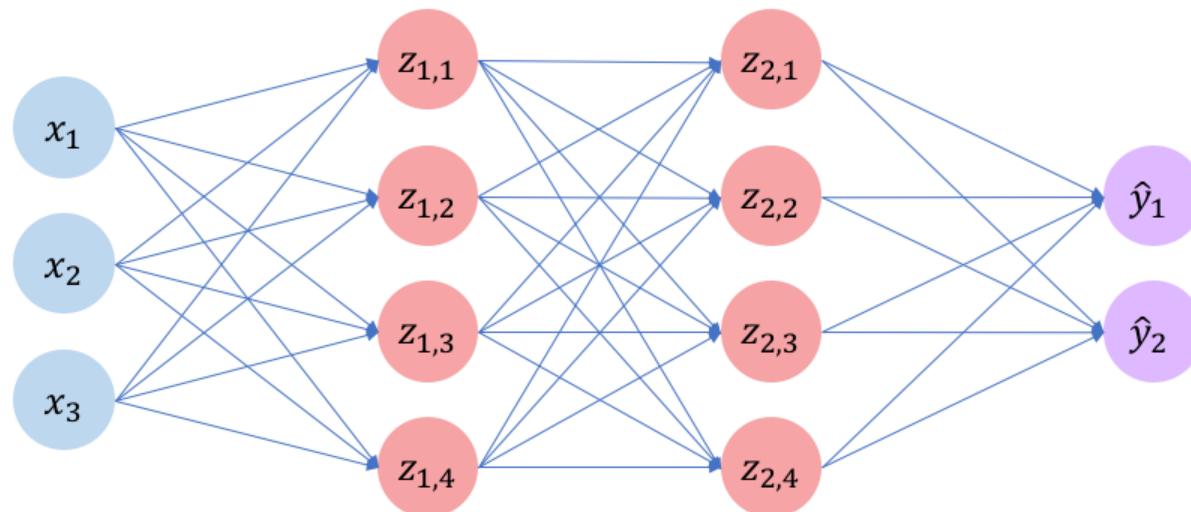
Technique that **constrains** our optimization problem to **discourage** complex models

Why we need it?

Improve **generalization** of our model on unseen data

Regularization 1 : Dropout

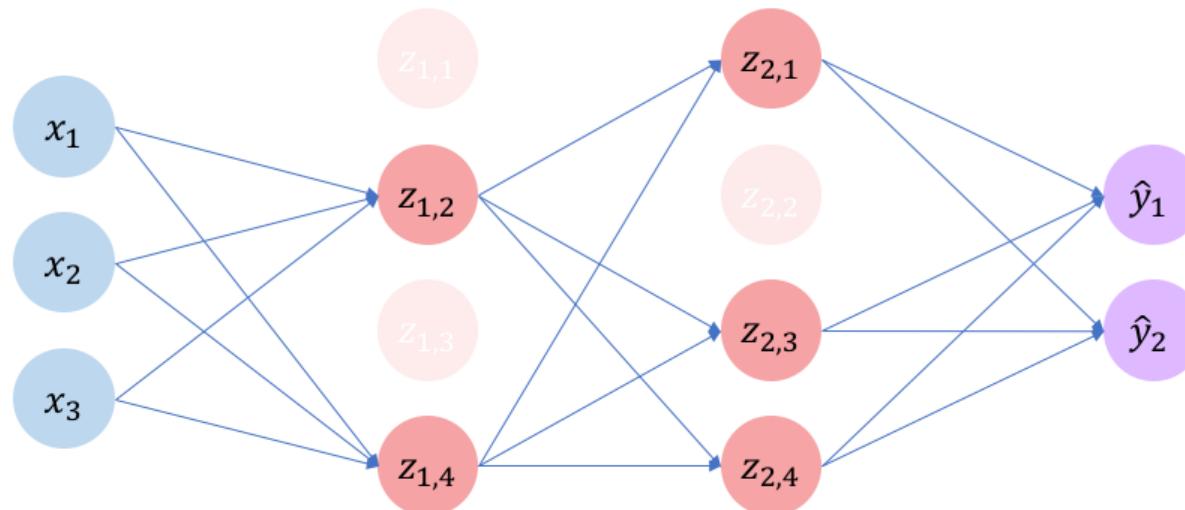
- During training, randomly set some activations to 0



Regularization 1 : Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

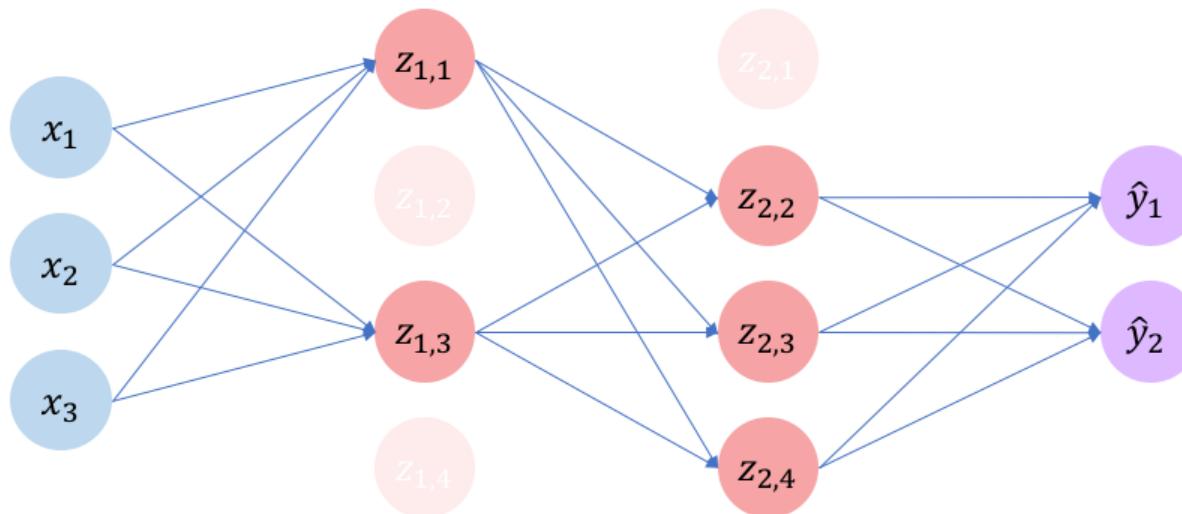
 `torch.nn.Dropout(p=0.5)`



Regularization 1 : Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

 `torch.nn.Dropout(p=0.5)`



Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



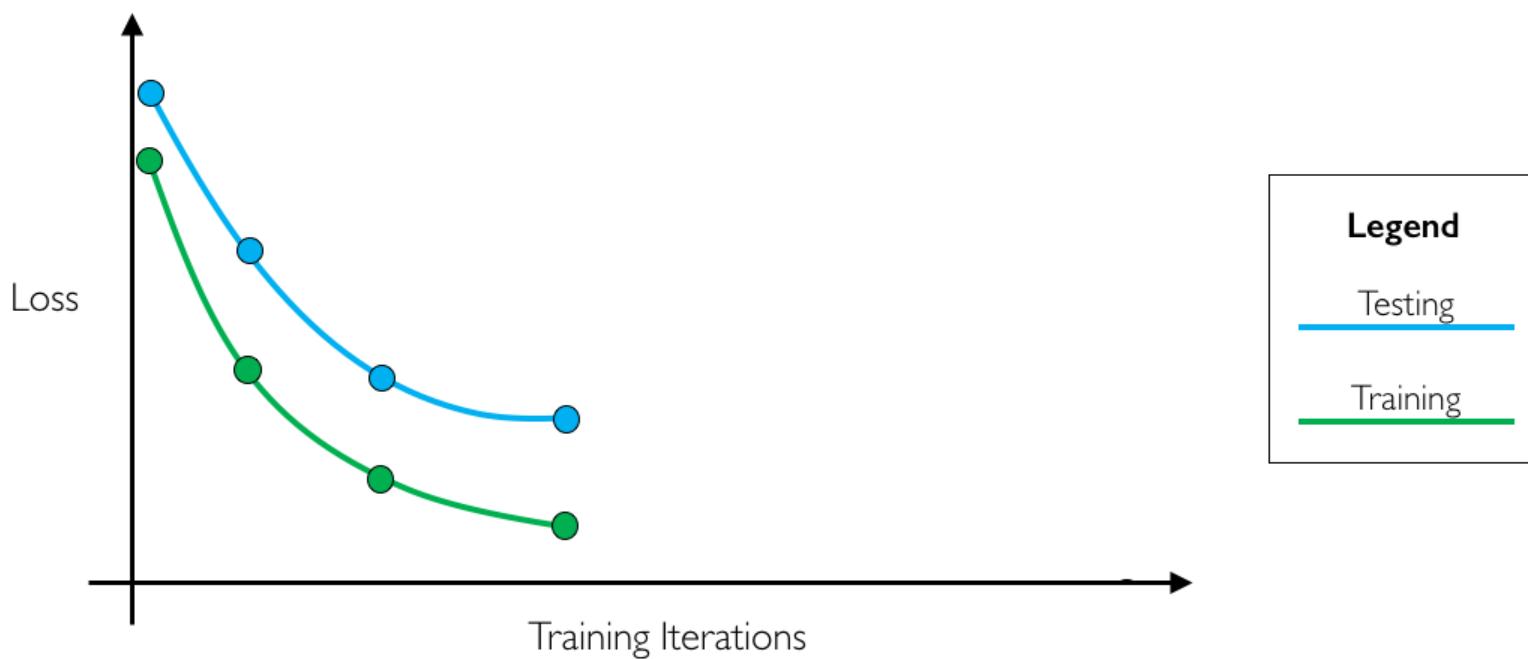
Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



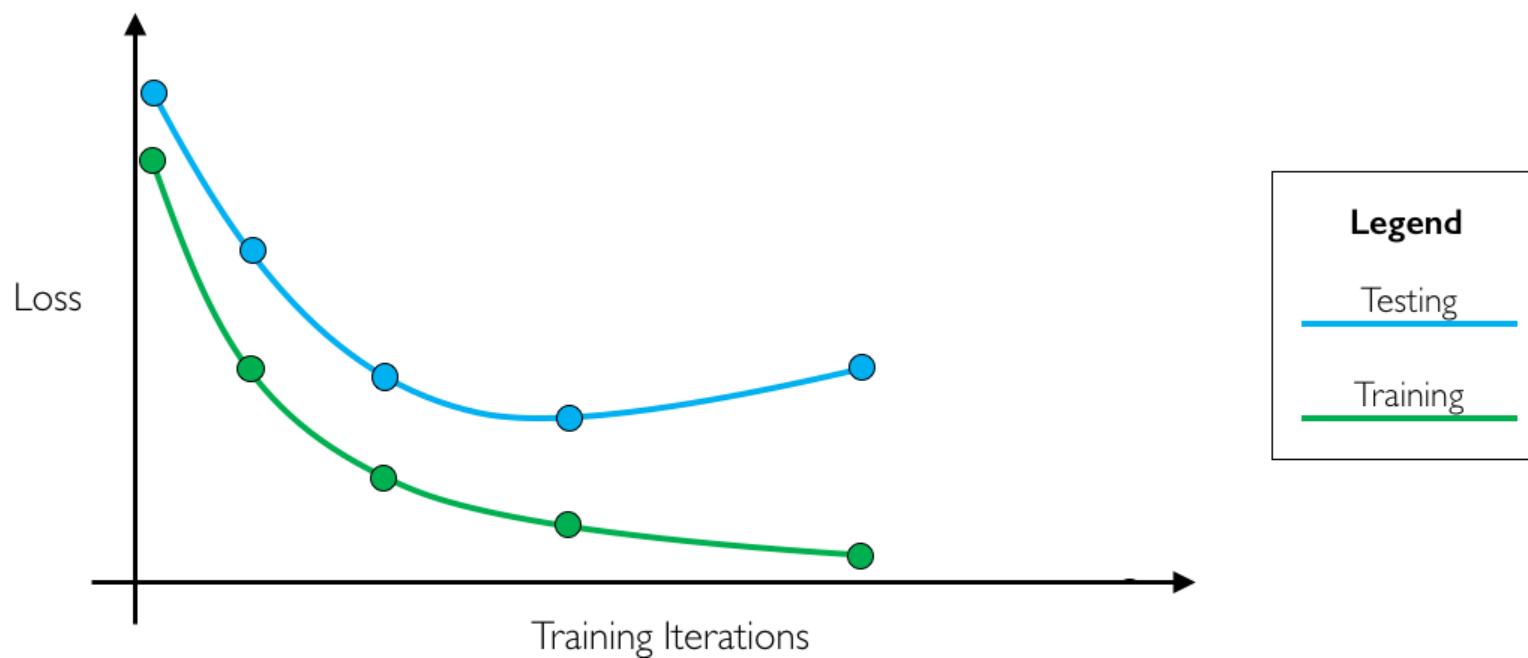
Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



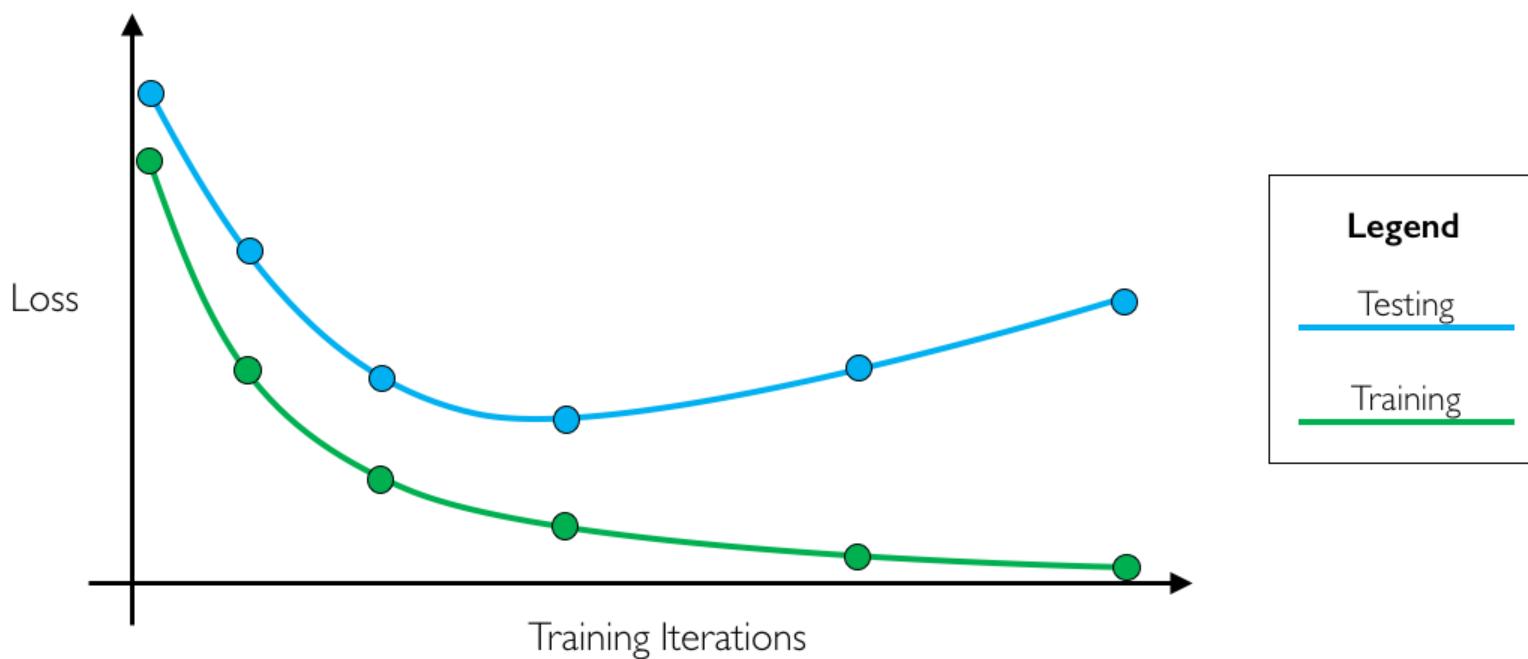
Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



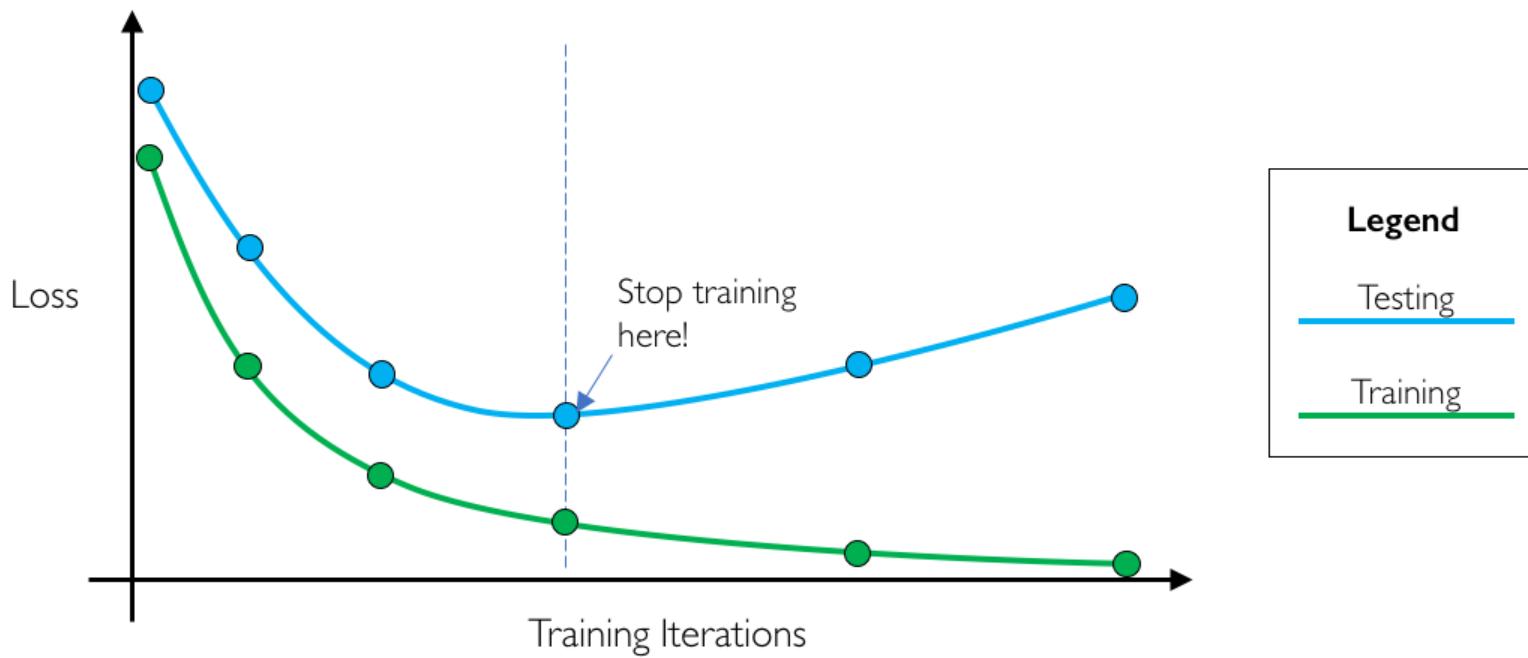
Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



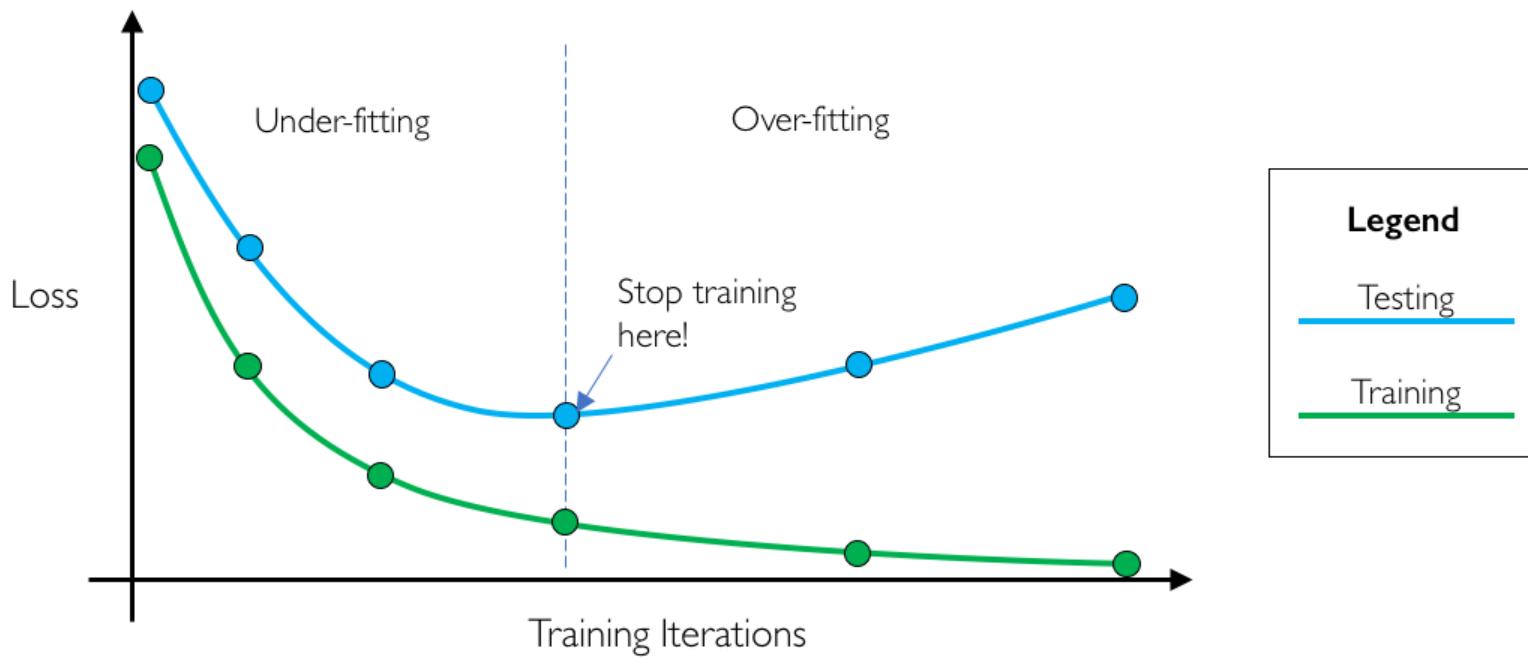
Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit

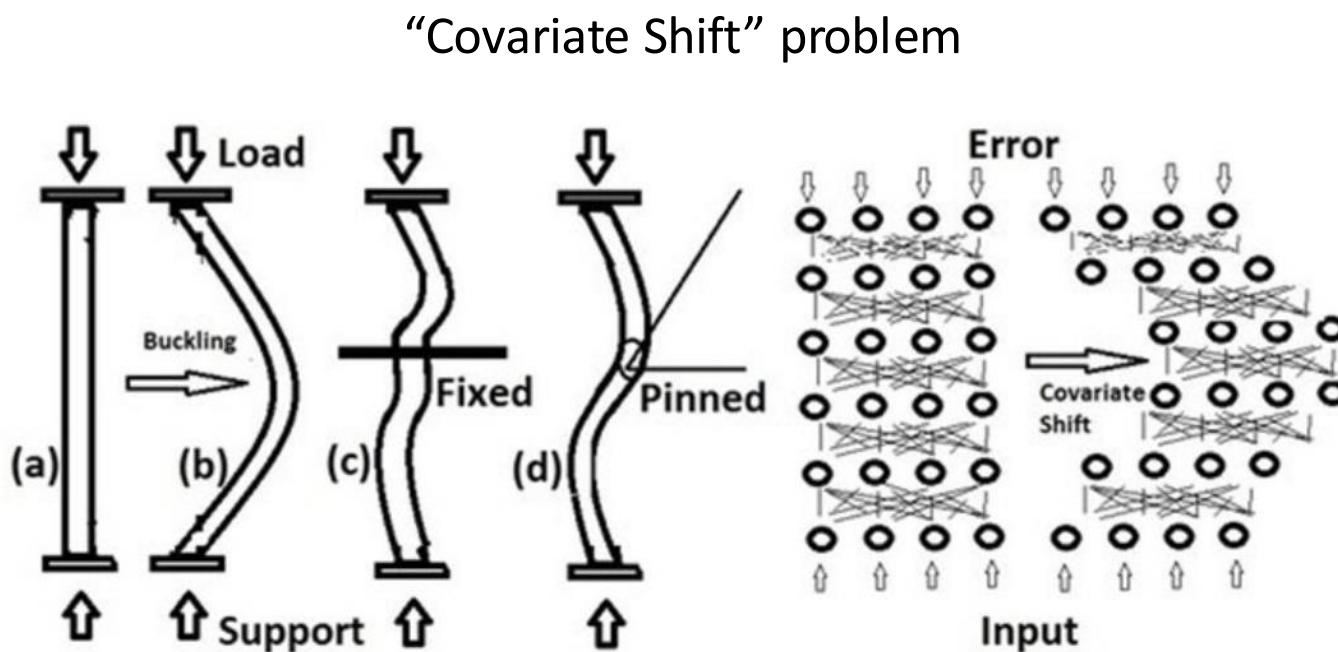


Regularization 2 : Early Stopping

- Stop training before we have a chance to overfit



Regularization 3 : Batch normalization



For both, Buckling or Co-Variate Shift a small perturbation leads to a large change in the later.

Debiprasad Ghosh, PhD, Uses AI in Mechanics

Regularization 3 : Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

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

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

 `torch.nn.BatchNorm2d(num_features=num_features)`

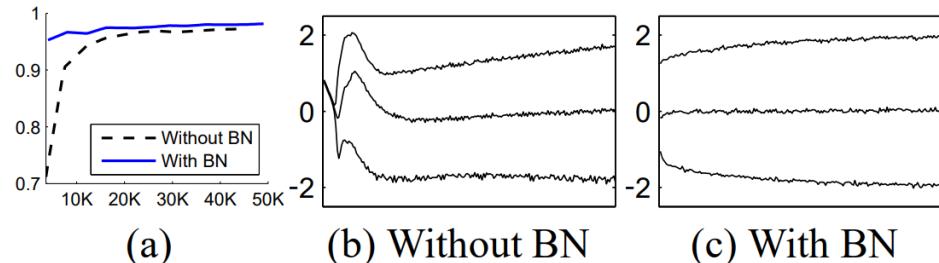
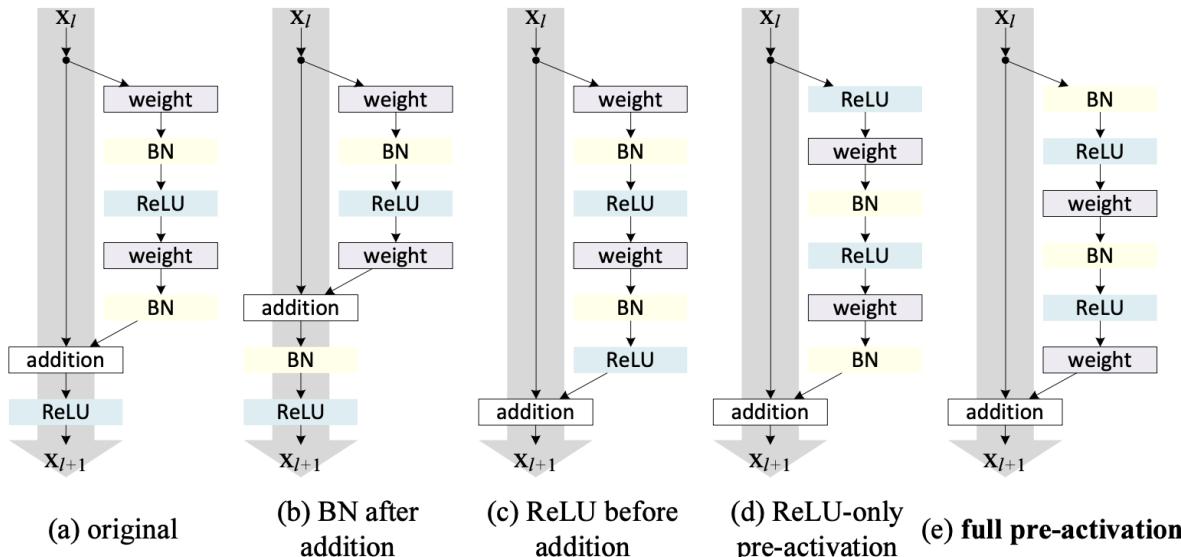


Figure 1: (a) The test accuracy of the MNIST network trained with and without Batch Normalization, vs. the number of training steps. Batch Normalization helps the network train faster and achieve higher accuracy. (b, c) The evolution of input distributions to a typical sigmoid, over the course of training, shown as {15, 50, 85}th percentiles. Batch Normalization makes the distribution more stable and reduces the internal covariate shift.

Regularization 3 : Batch normalization

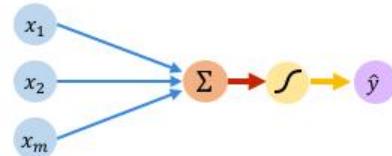


```
conv1 = torch.nn.Sequential{  
    torch.nn.Conv2d(...),  
    torch.nn.BatchNorm2d(num_features),  
    torch.nn.relu(...)  
}
```

Core Foundation Review

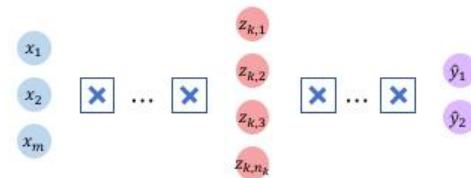
The Perceptron

- Structural building blocks
- Nonlinear activation functions



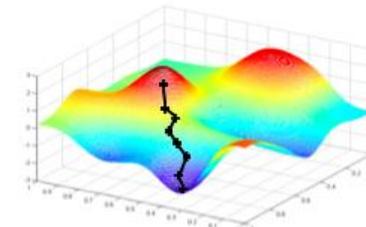
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization



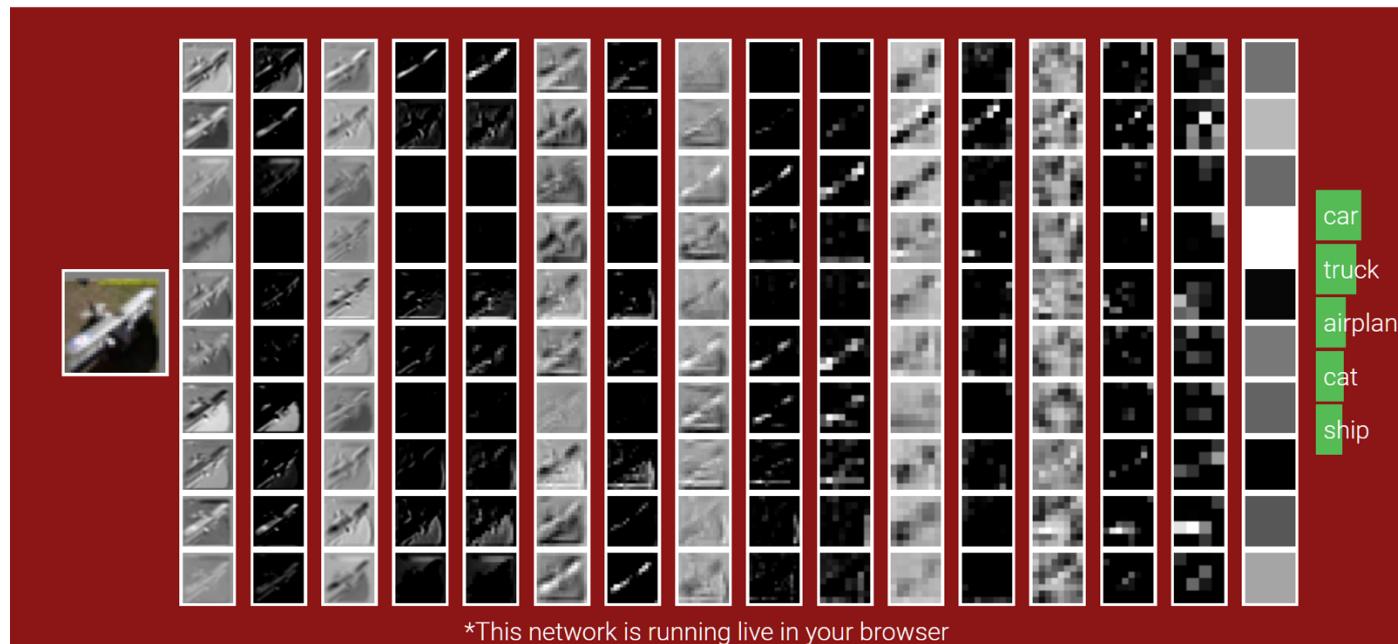
추천 학습 자료

- Stanford 231n by Fei-Fei Li

CS231n: Convolutional Neural Networks for Visual Recognition

Spring 2020

Previous Years: [\[Winter 2015\]](#) [\[Winter 2016\]](#) [\[Spring 2017\]](#) [\[Spring 2018\]](#) [\[Spring 2019\]](#)



<http://cs231n.stanford.edu/index.html>

추천 학습 자료

- Stanford cs224n by Chris Manning

[CS224n Home](#)

Coursework

Schedule

Office Hours

Lecture Videos

Project Reports

Piazza

CS224n: Natural Language Processing with Deep Learning

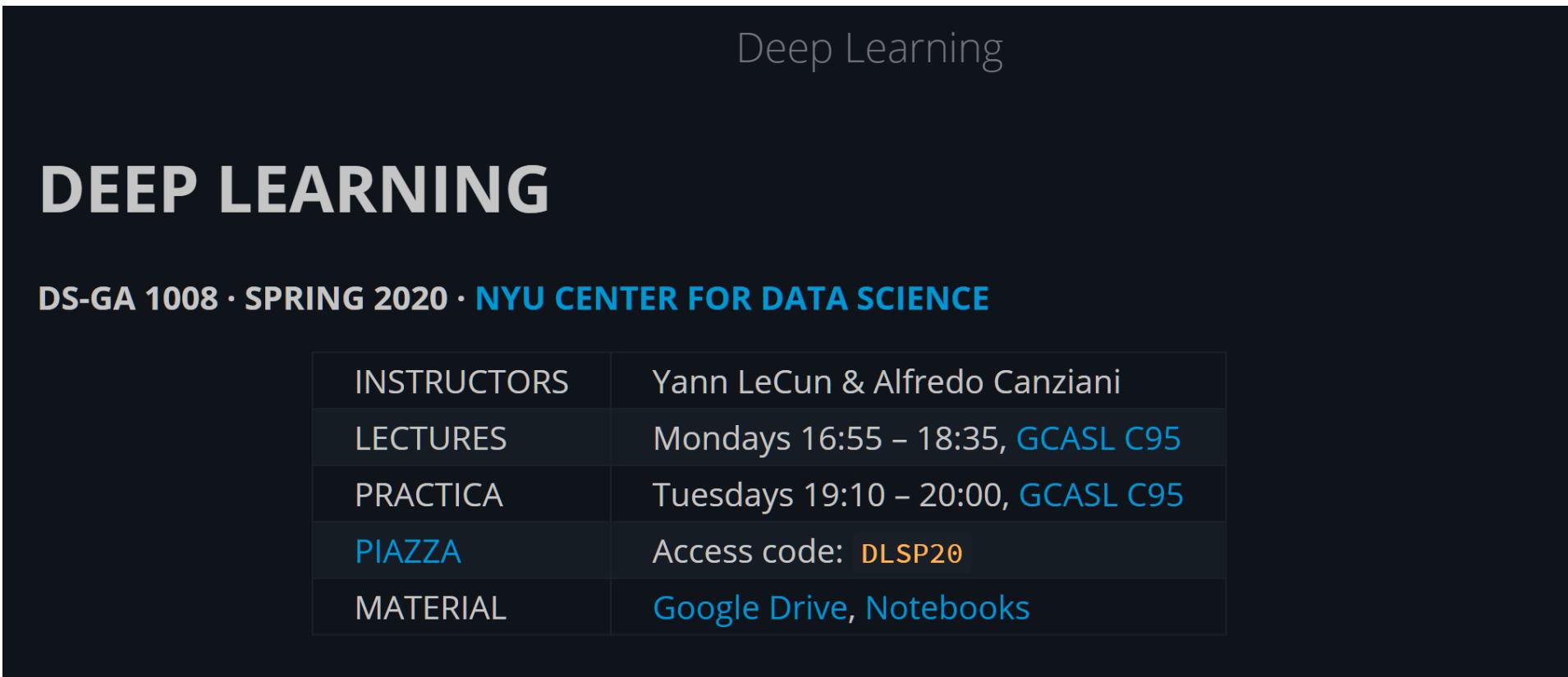
Stanford / Winter 2020

Natural language processing (NLP) is a crucial part of artificial intelligence (AI), modeling how people share information. In recent years, deep learning approaches have obtained very high performance on many NLP tasks. In this course, students gain a thorough introduction to cutting-edge neural networks for NLP.

<http://web.stanford.edu/class/cs224n/index.html#schedule>

추천 학습 자료

- NYU DEEP LEARNING by Yann Lecun



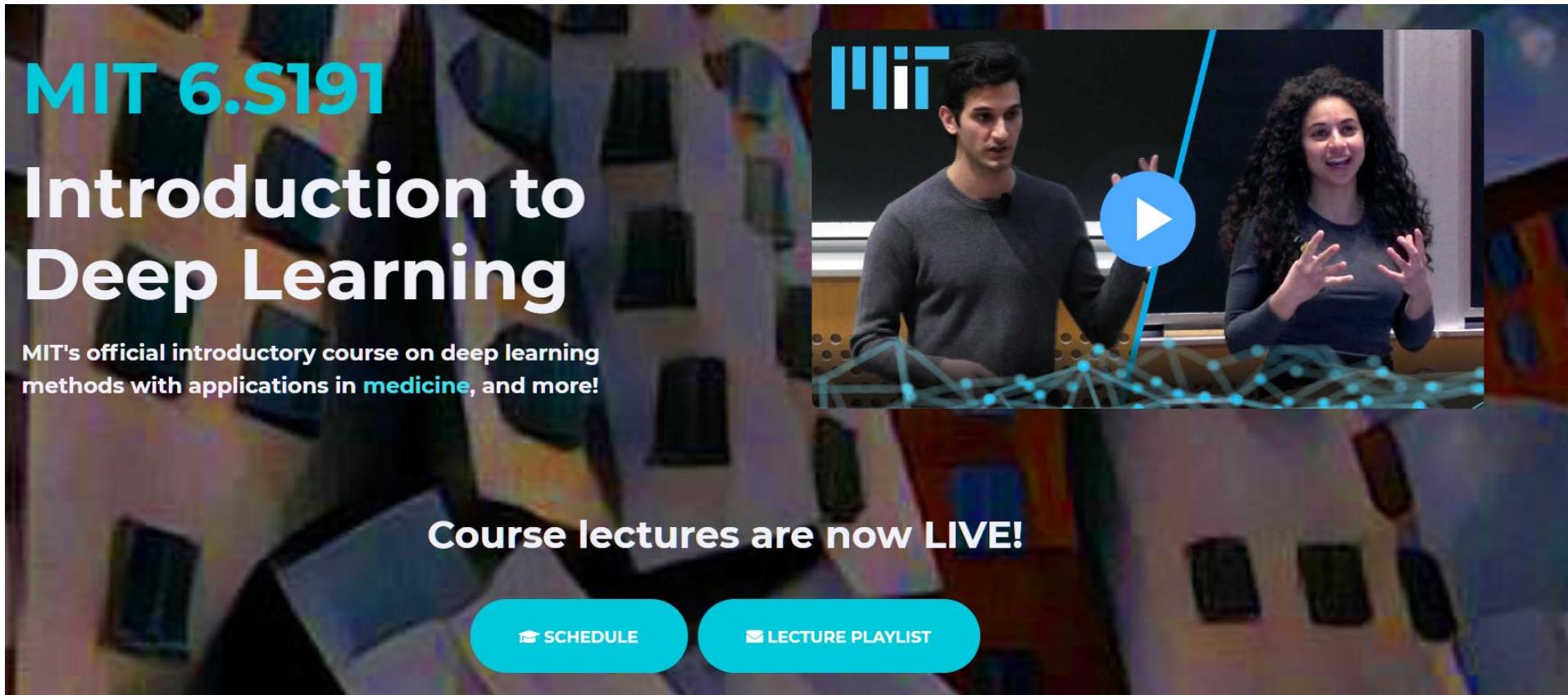
The screenshot shows the homepage of the NYU Deep Learning course. At the top right, it says "Deep Learning". The main title "DEEP LEARNING" is in large white letters. Below it, the course information "DS-GA 1008 · SPRING 2020 · NYU CENTER FOR DATA SCIENCE" is displayed. A table provides details about the course structure:

INSTRUCTORS	Yann LeCun & Alfredo Canziani
LECTURES	Mondays 16:55 – 18:35, GCASL C95
PRACTICA	Tuesdays 19:10 – 20:00, GCASL C95
PIAZZA	Access code: DLSP20
MATERIAL	Google Drive , Notebooks

<https://atcold.github.io/pytorch-Deep-Learning/>

추천 학습 자료

- MIT



추천 학습 자료

- <https://deep-learning-drizzle.github.io/index.html#dldnn>



Deep Learning Drizzle



"Read enough so you start developing intuitions and then trust your intuitions and go for it!" 

Prof. Geoffrey Hinton, University of Toronto

Deep Learning (Deep Neural Networks) 	Probabilistic Graphical Models 
Machine Learning Fundamentals 	Natural Language Processing 
Optimization for Machine Learning 	Automatic Speech Recognition 
General Machine Learning 	Modern Computer Vision 
Reinforcement Learning 	Boot Camps or Summer Schools 
Bayesian Deep Learning 	Medical Imaging 
Graph Neural Networks 	Bird's-eye view of Artificial Intelligence 

References

- <http://cs231n.stanford.edu/index.html>
- <http://web.stanford.edu/class/cs224n/index.html#schedule>
- <https://atcold.github.io/pytorch-Deep-Learning/>
- http://introtodeeplearning.com/?fbclid=IwAR1_iEF374Jg46s_2HMVWcVPFGwZWorBxEL25jzfpFfgcIFWbo-BDWUv_54