

# **Intro to Deep Learning**

Sungwook E. Hong (KASI)

2025. 8. 19 @ A3Net Summer School

# What is Learning?

Supervised

From given sets of  $\{\vec{x}_i\}$  and  $\{\vec{y}_i\}$ ,  
find a function  $F(\dots)$  that satisfies

$$F(\vec{x}_i) \approx \vec{y}_i$$

Outputs that we  
need to reproduce

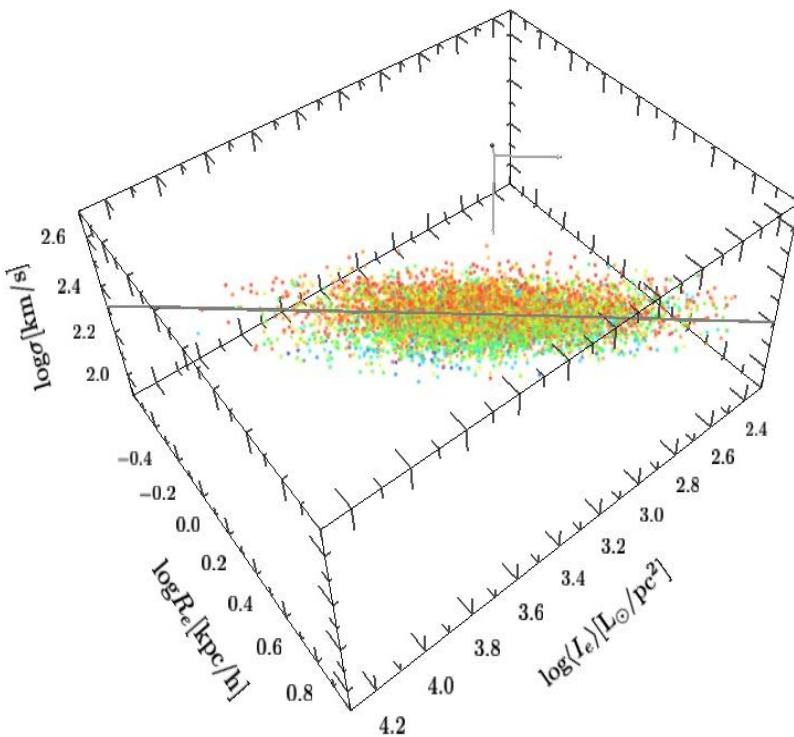
Unsupervised

Inputs

From a given set of  $\{\vec{x}_i\}$ ,  
find a function  $F(\dots)$  that produces a useful  $\{\vec{y}_i\}$  by

$$F(\vec{x}_i) = \vec{y}_i$$

Outputs that  
we can define

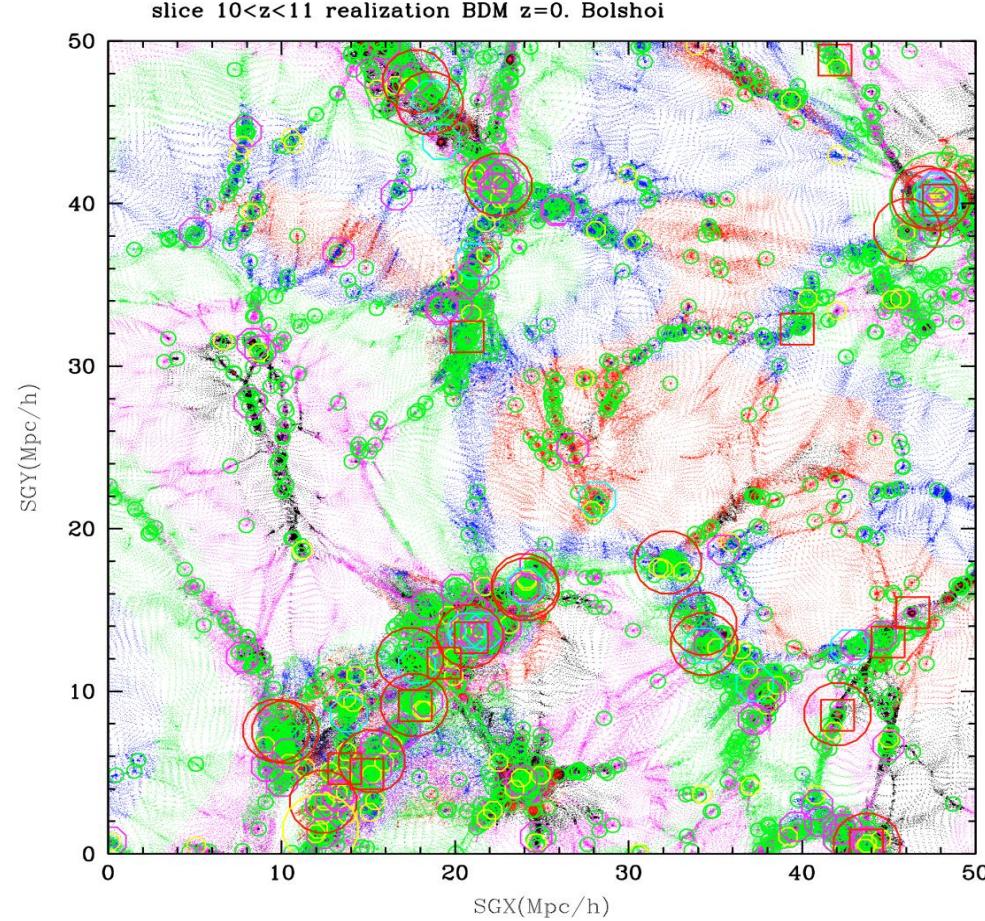


Fitting with a given functional form

→ **Supervised**

**Inputs:** effective radius & velocity dispersion

**Output:** surface brightness



Finding dark matter halos

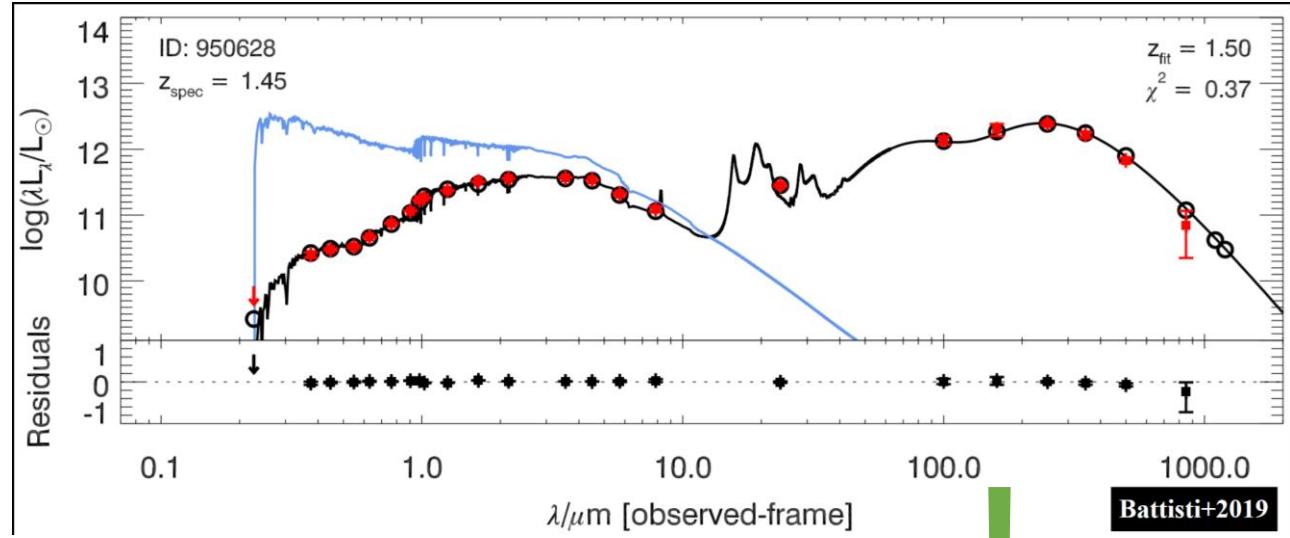
→ **Unsupervised**

**Inputs:** particle positions

**Outputs:** halo IDs

# In Classical (Machine) Learning...

Wavelength	Intensity
X.XXXXXXX	X.XXXXXXX
...	...



Find fitting formulae based on theory & experiences

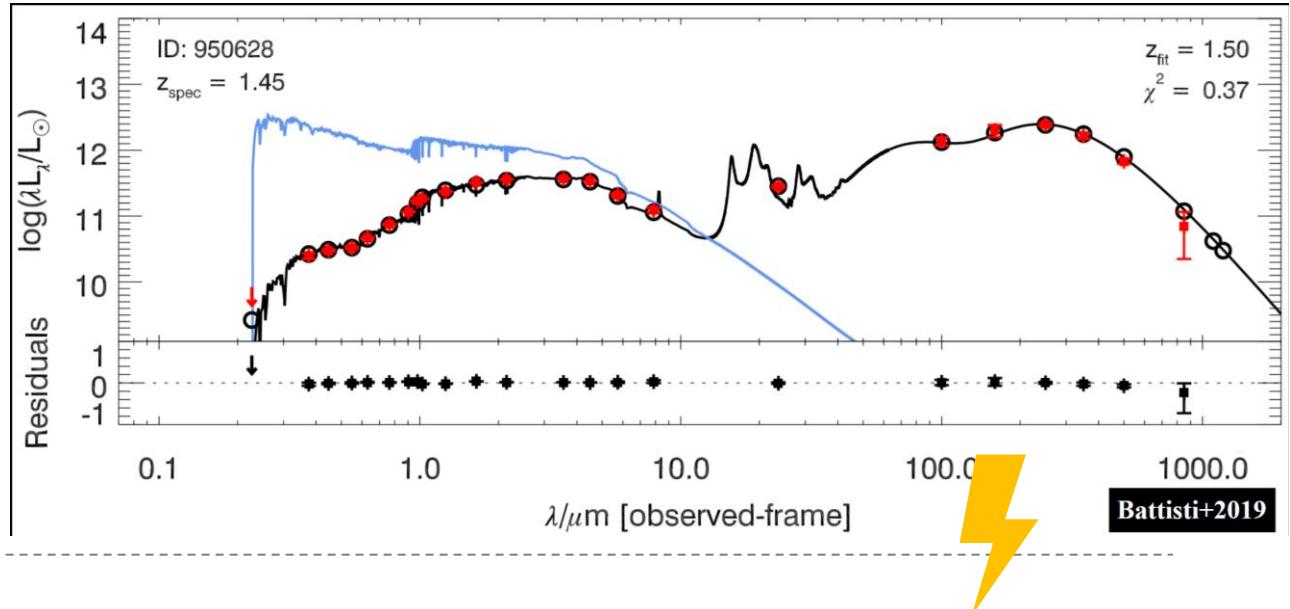


Object type / redshift / ...



# In Machine Learning...

Wavelength	Intensity
X.XXXXXXX	X.XXXXXXX
...	...



Set the learning strategies



Object type / redshift / ...

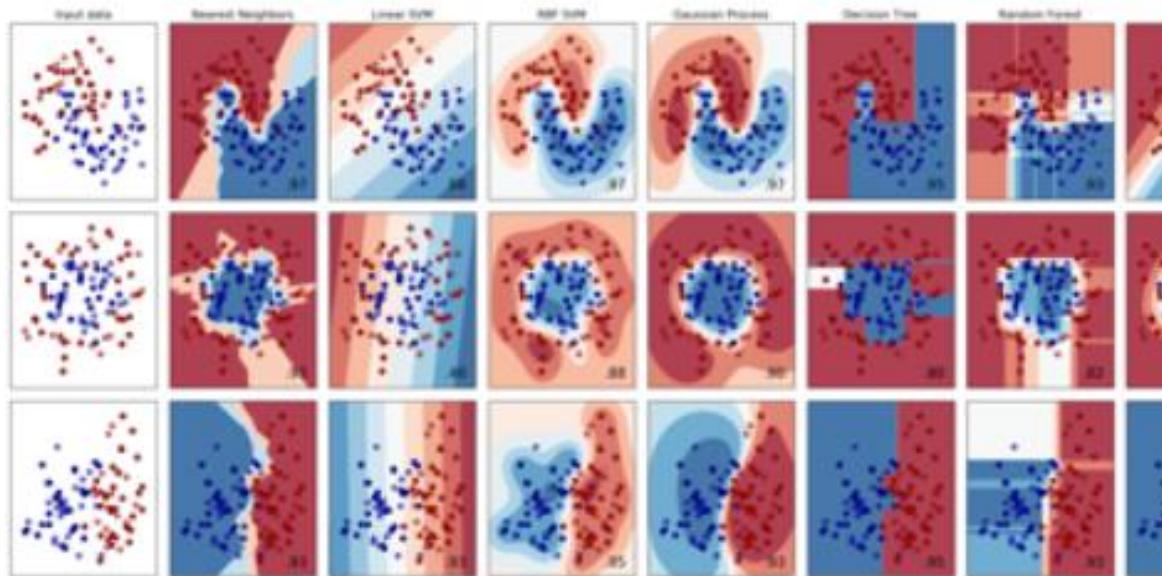


# Classification

Identifying which category an object belongs to.

**Applications:** Spam detection, image recognition.

**Algorithms:** SVM, nearest neighbors, random forest, and more...



**Type: Supervised**

**Outputs: Vector**

Each element means the probability of being each type

**Example**

Classify each image to either **red**, **green**, or **blue**!

**Inputs**



**Outputs**

(1, 0, 0)



(0, 1, 0)



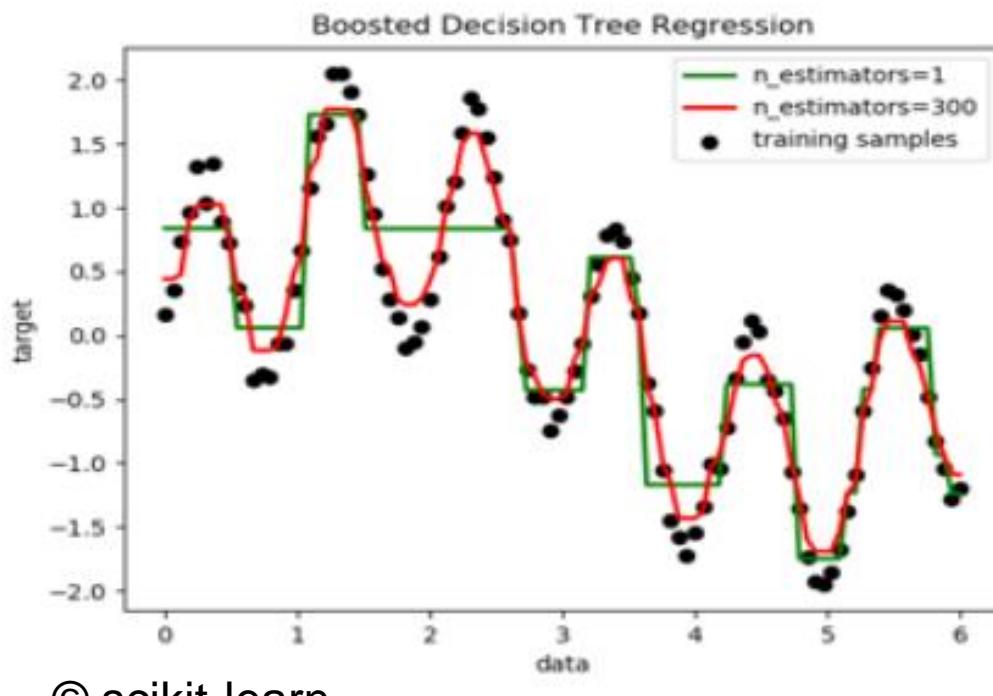
(0, 0, 1)

# Regression

Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.

**Algorithms:** SVR, nearest neighbors, random forest, and more...

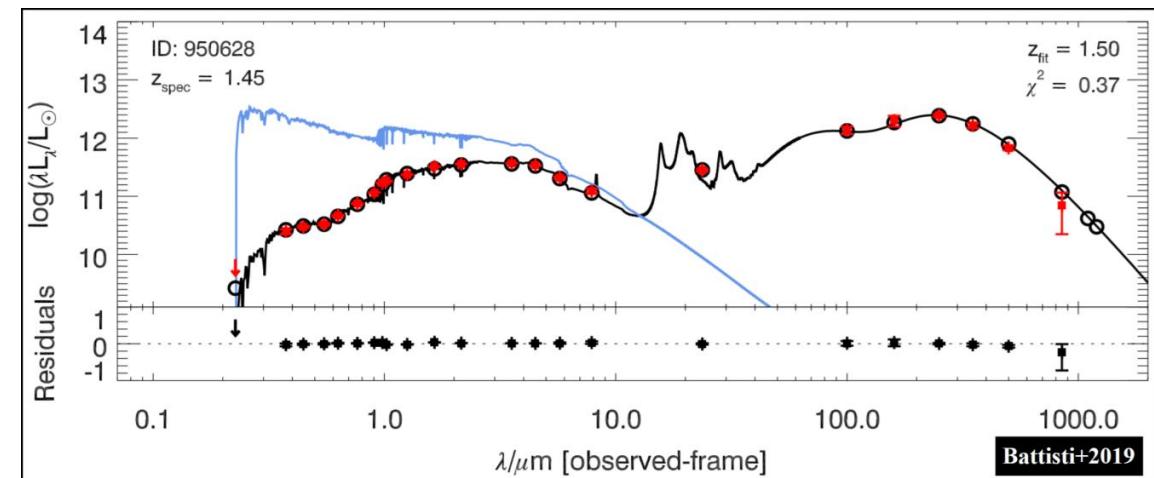


Type: Supervised

Outputs: Vector

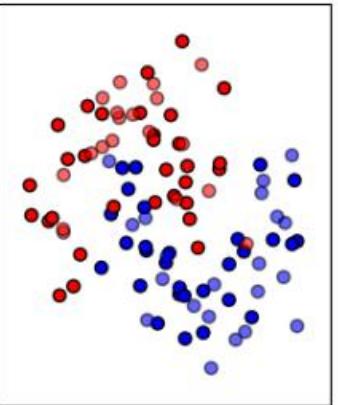
Example

Find the object type and redshift from the spectrum

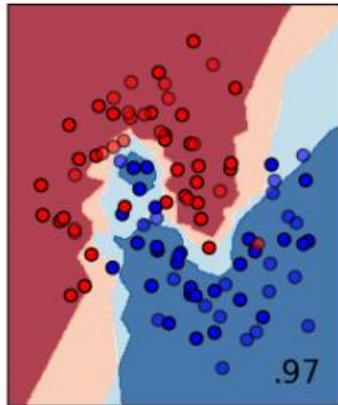


Classification  $\longleftrightarrow$  Regression

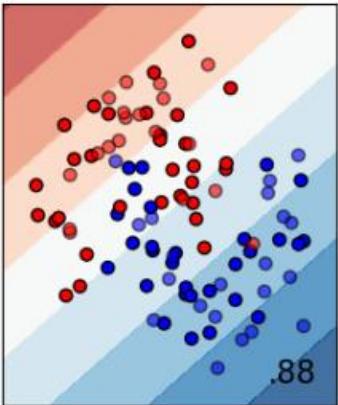
Input data



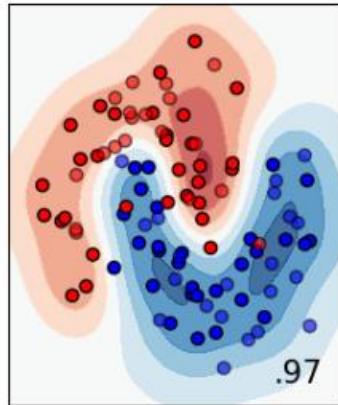
Nearest Neighbors



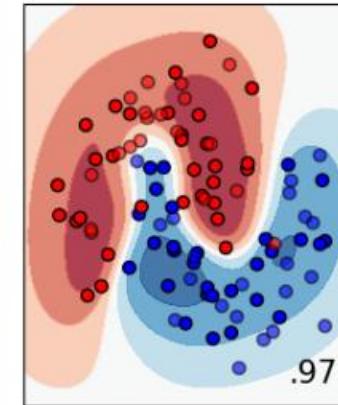
Linear SVM



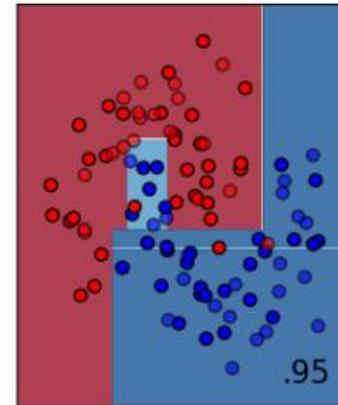
RBF SVM



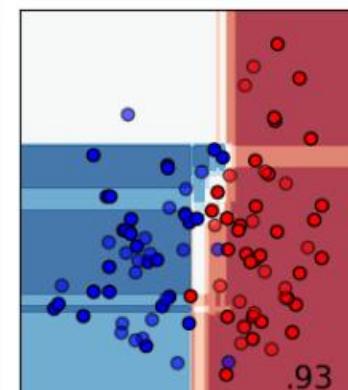
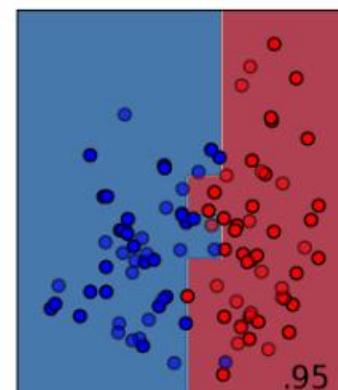
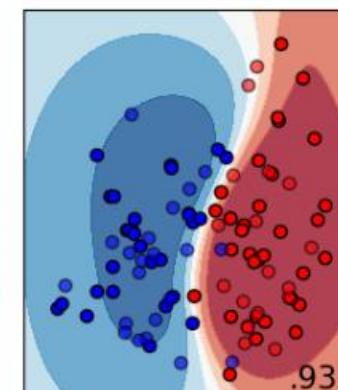
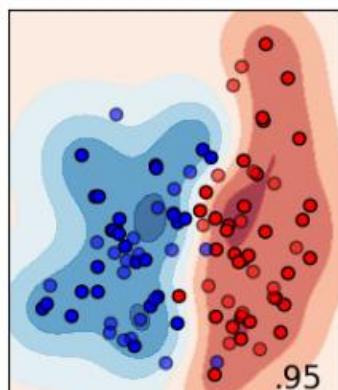
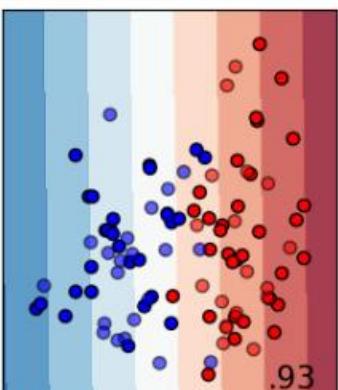
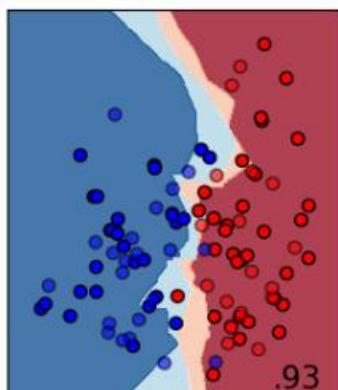
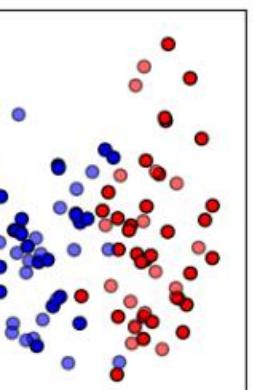
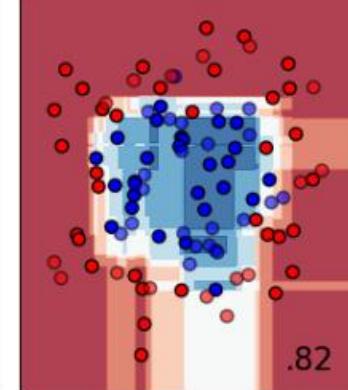
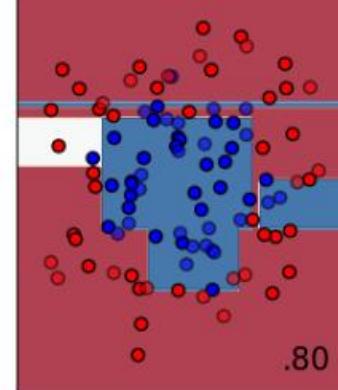
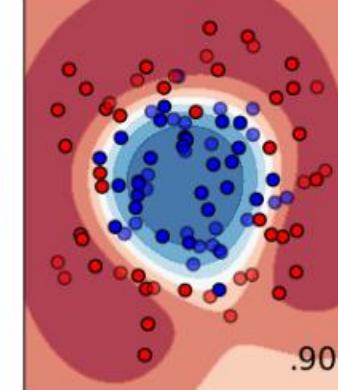
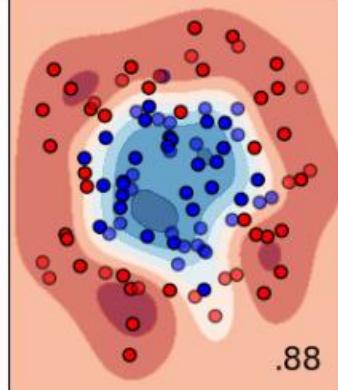
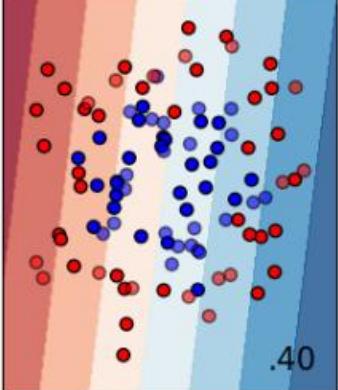
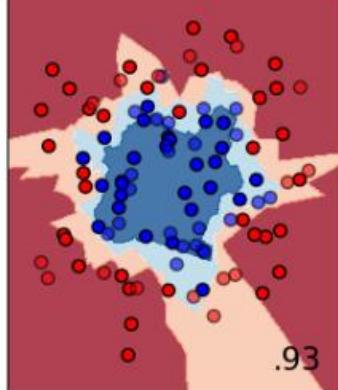
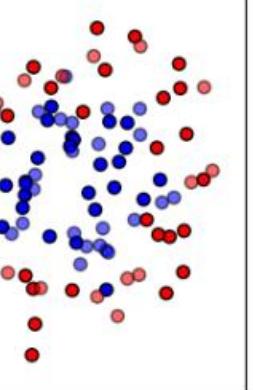
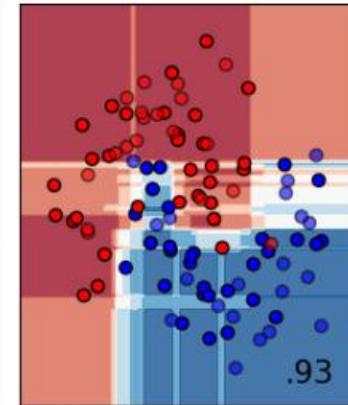
Gaussian Process



Decision Tree



Random Forest



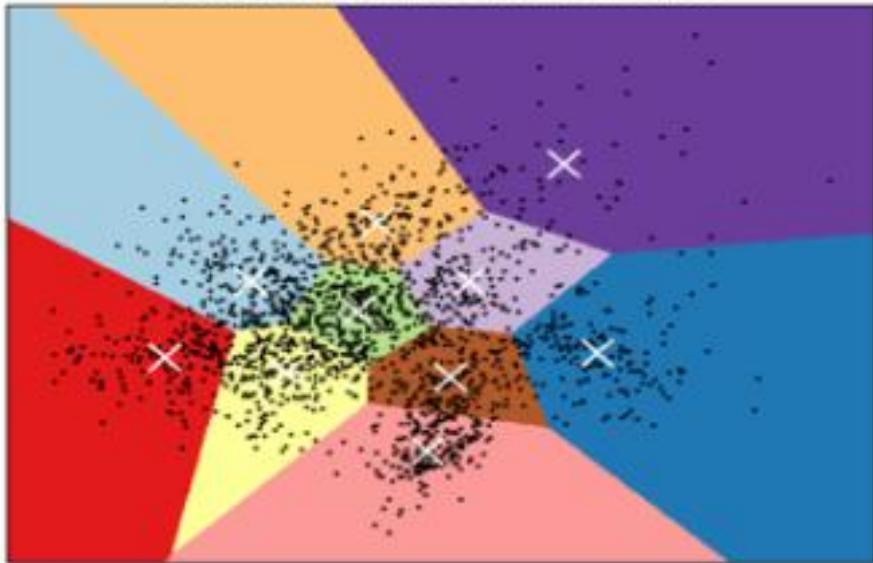
## Clustering

Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes

**Algorithms:** k-Means, spectral clustering, mean-shift, and more...

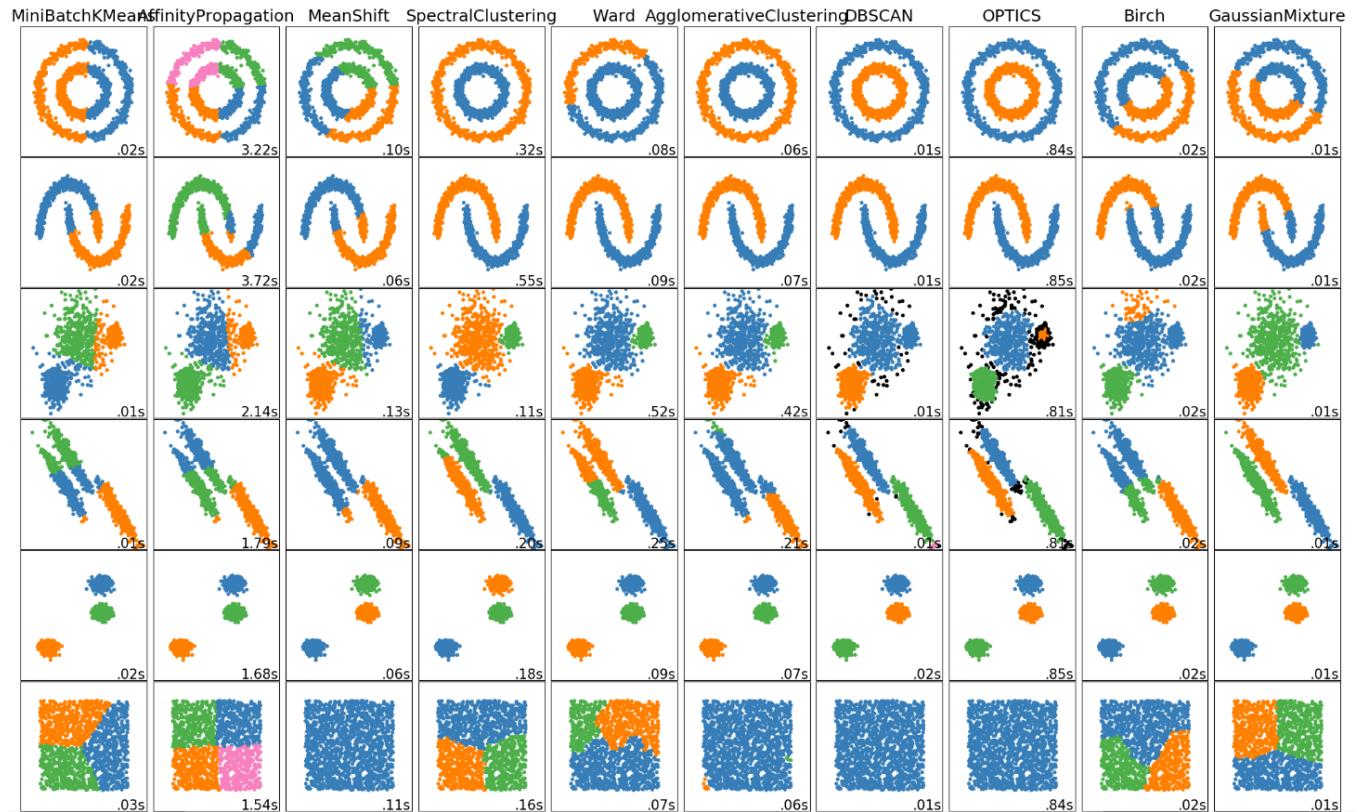
K-means clustering on the digits dataset (PCA-reduced data)  
Centroids are marked with white cross



© scikit-learn

## Type: Unsupervised

# Outputs: Group ID

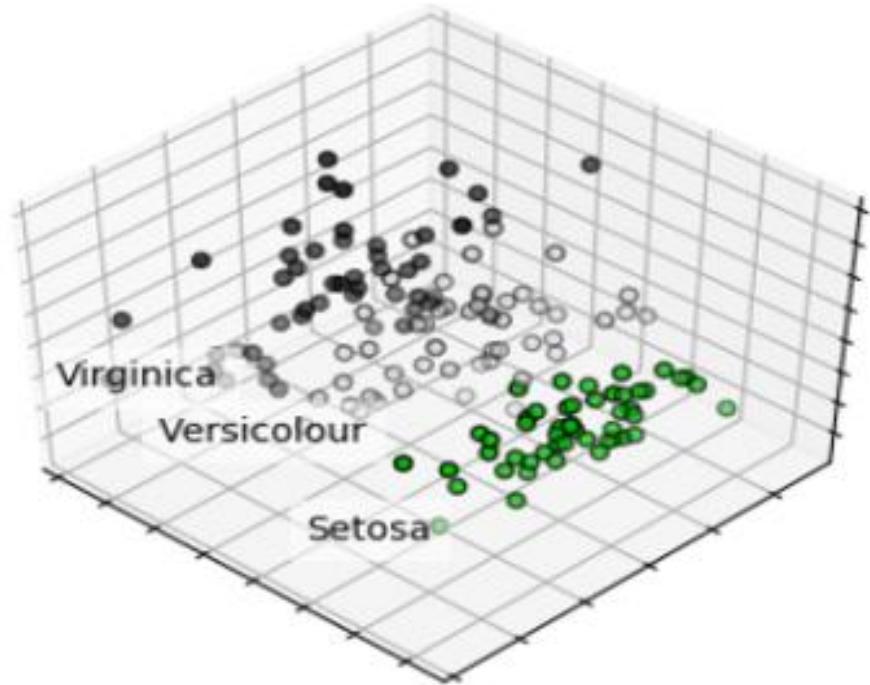


# Dimensionality reduction

Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency

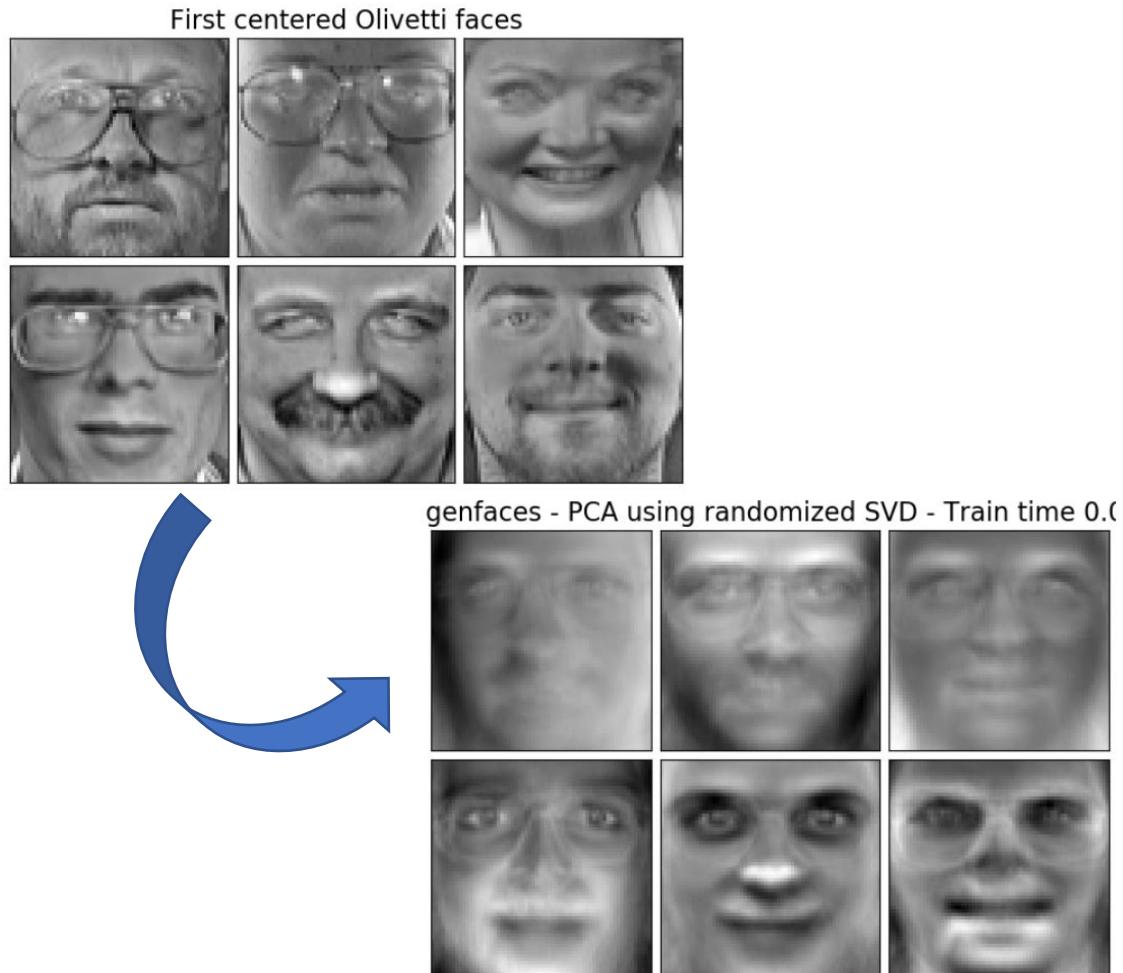
**Algorithms:** k-Means, feature selection, non-negative matrix factorization, and more...



© scikit-learn

# Type: Unsupervised

## Outputs: Vectors with lower dimensions



# Machine Learning

- All results should be lead by physics
  - Software is no replacement for physical understanding
- ML is a useful technique in three cases:
  - No known or understandable physical mechanism exists:
    - Dark energy
  - Physics known, but computing the required quantities is challenging:
    - Photometric redshifts
    - Large numbers of simulations - non-linear emulators
  - Paradigm testing: looking for things not predicted by any physics or models
    - Unknown unknown
    - Model-independent tests

Techniques (examples):

## 1. Classification/catagoriszation

- Support Vector Machines
- Naive Bayes
- Perceptron/CNN
- Decision Trees

## 2. Clustering

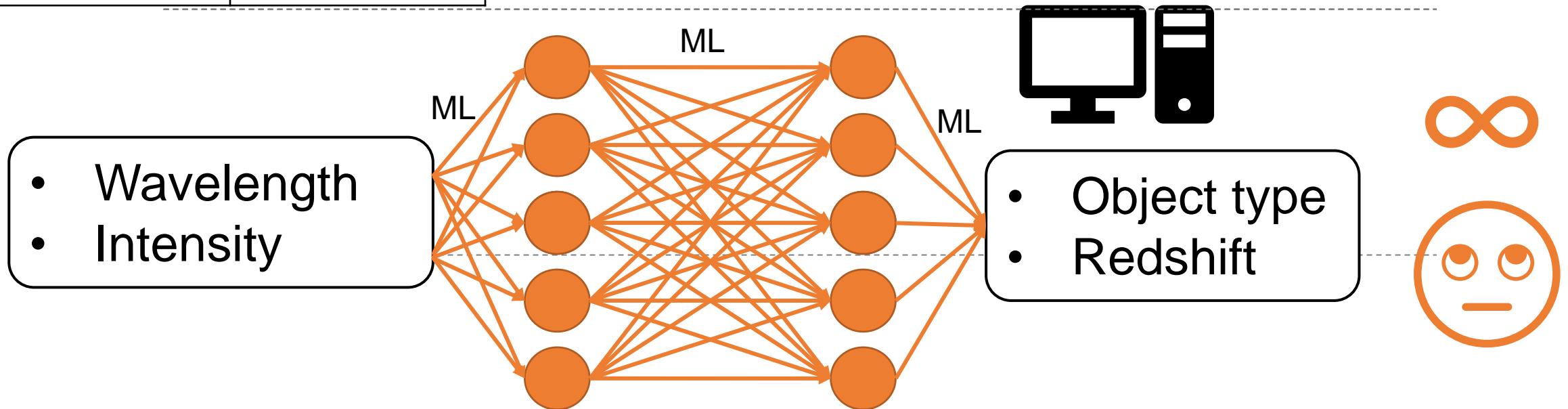
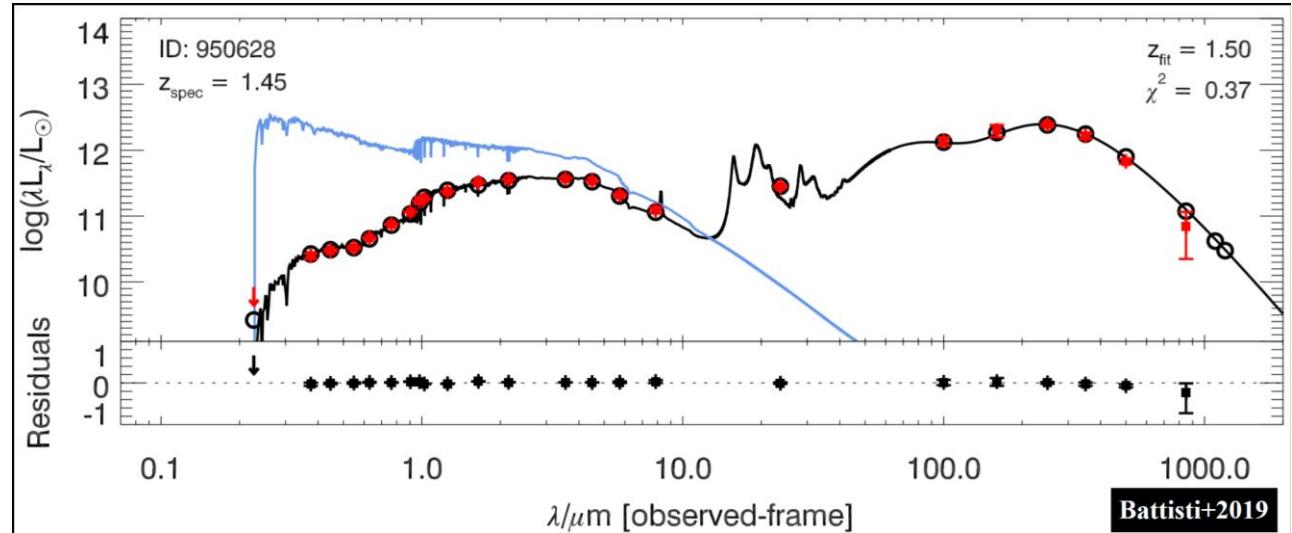
- Hierarchical
- k-means
- DBSCAN
- Density

## 3. Interpolation/Regression

- Gaussian Process Regression
- Kernel Ridge Regression
- Multivariate adaptive regression splines

# Then, In Deep Learning...

Wavelength	Intensity
X.XXXXXXX	X.XXXXXXX
...	...



# Compared to ML, DL is...

- Harder to understand the underlying physics
- Having too many degrees of freedom to build models
- Requiring more data and computational resources (> GPUs) for successful learning
- Harder to guarantee the success of learning mathematically
- Capable of solving very difficult problems, once successful

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

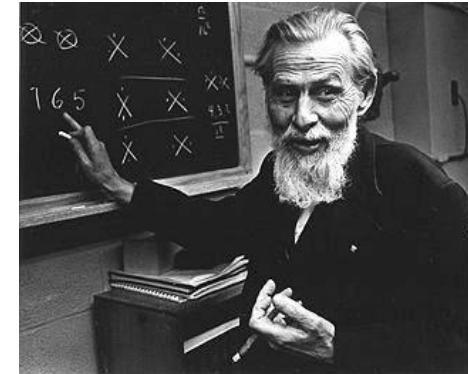
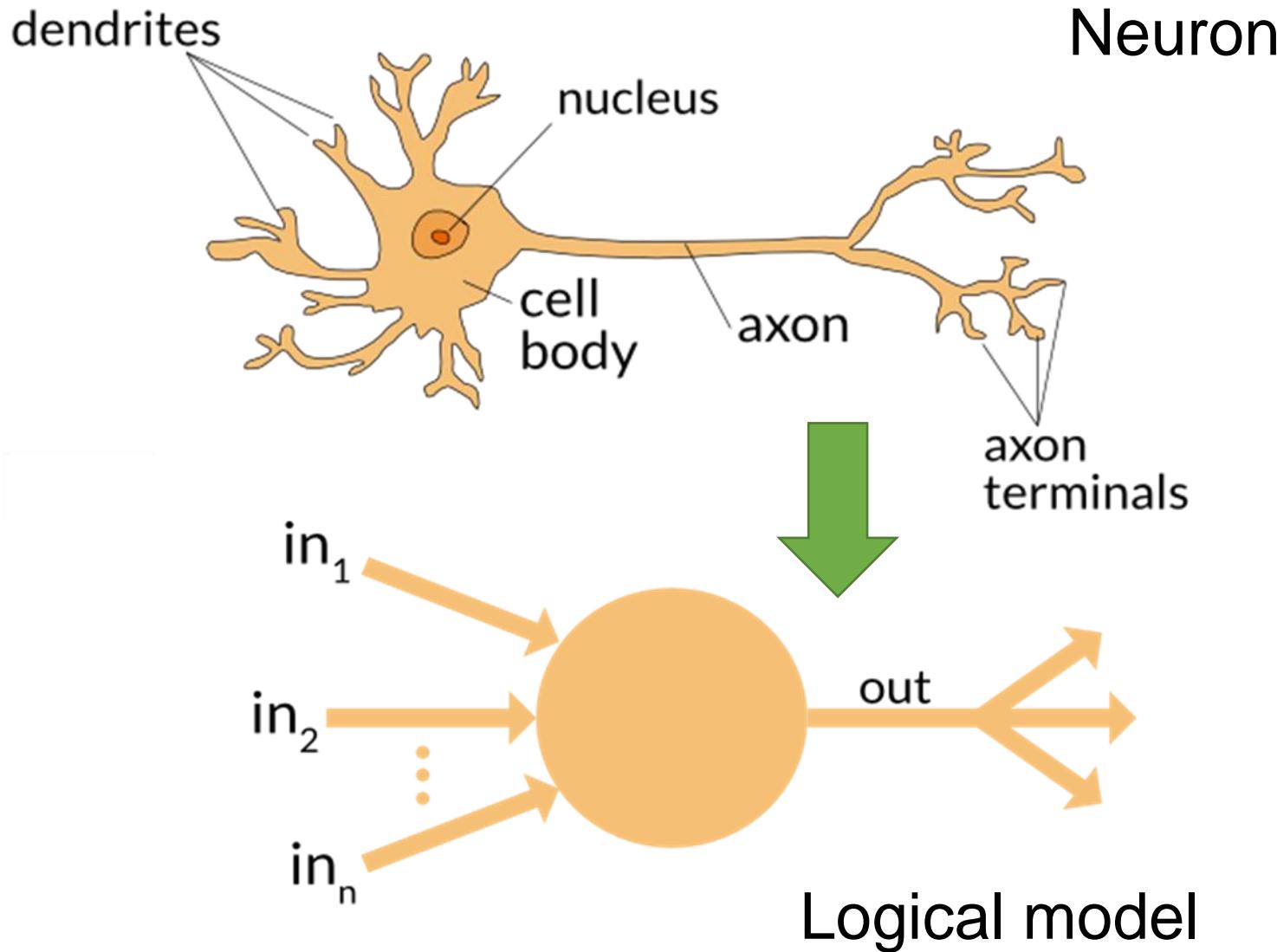
WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.



So...  
What is  
deep  
learning?

# A brief history... in 1943

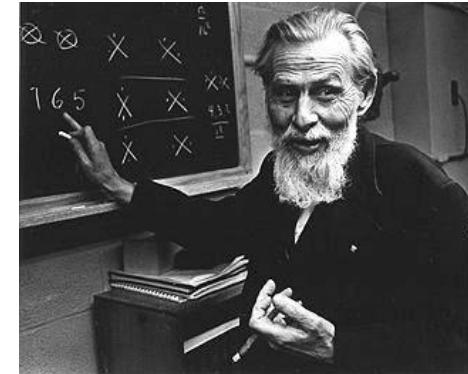
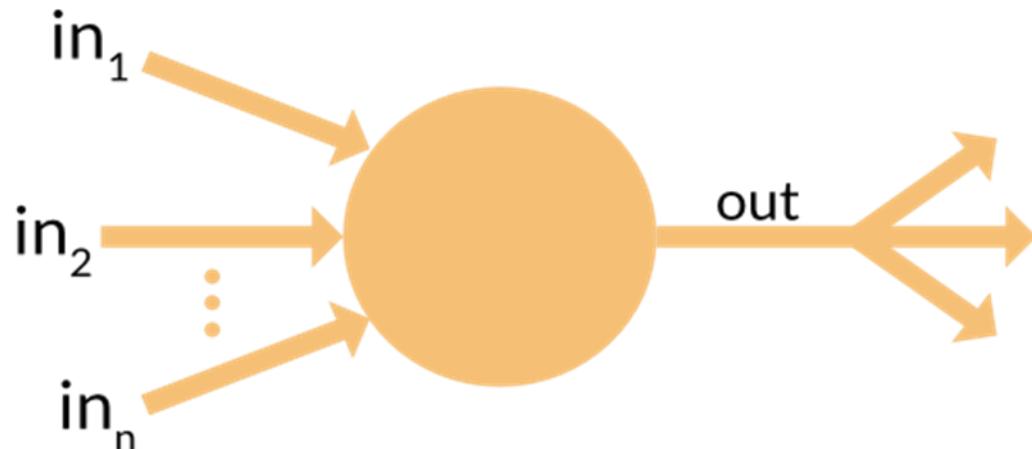


Warren McCulloch



Walter Pitts

# A brief history... in 1943



Warren McCulloch

**Activation**

$$Y = \begin{cases} 1, & \sum_{i \in E} X_i \geq \Theta \text{ and } \sum_{j \in I} X_j = 0 \\ 0, & \sum_{i \in E} X_i < \Theta \text{ or } \sum_{j \in I} X_j > 0 \end{cases}$$

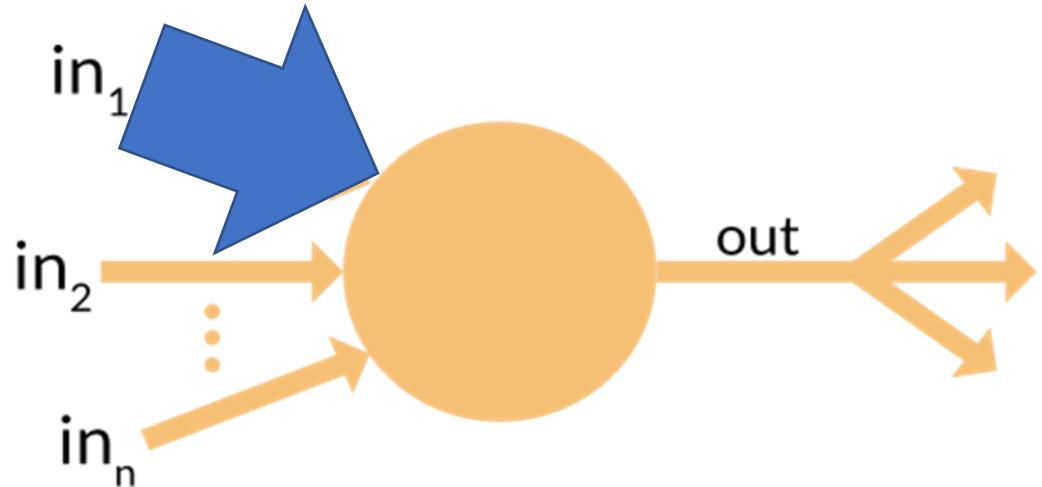
**Threshold**

**Excitation**      **Inhibition**



Walter Pitts

# A brief history... in 1949



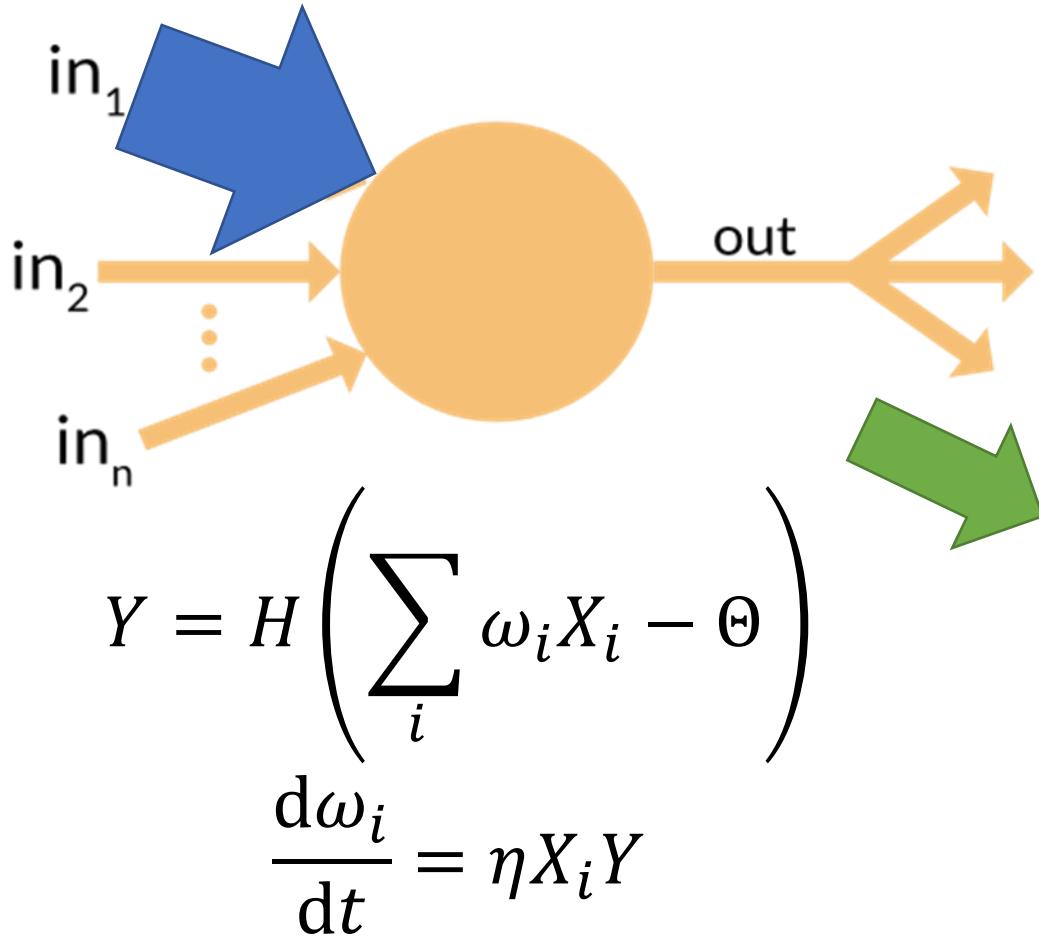
Donald Hebb

$$Y = H \left( \sum_i \omega_i X_i - \Theta \right)$$

A graph showing a green step function that increases at a specific point, representing the change in synaptic weight over time according to Hebb's rule.

$X_i \uparrow$  and  $Y \uparrow \Rightarrow \frac{d\omega_i}{dt} > 0$  (Synaptic strength)

# A brief history... in 1957

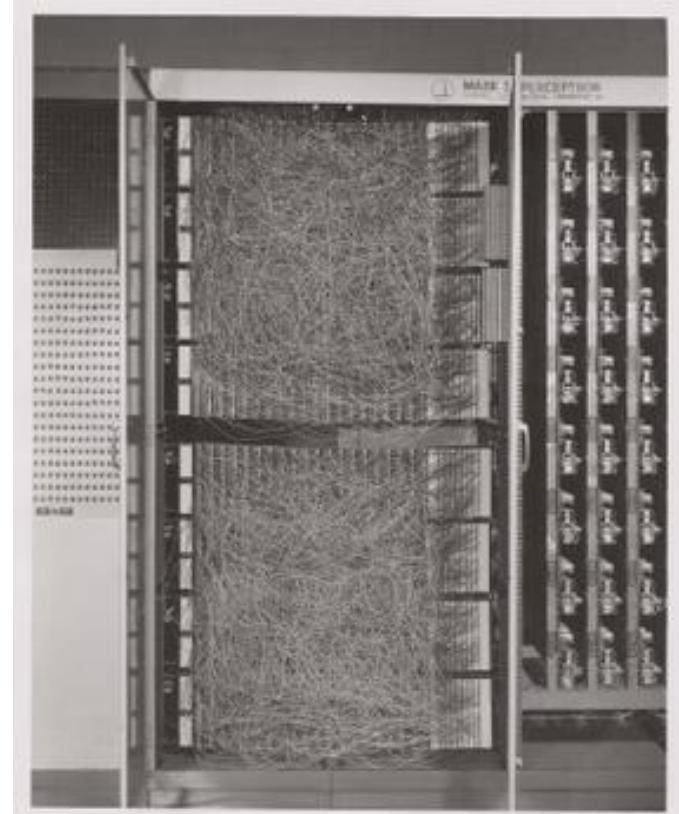


## Perceptron

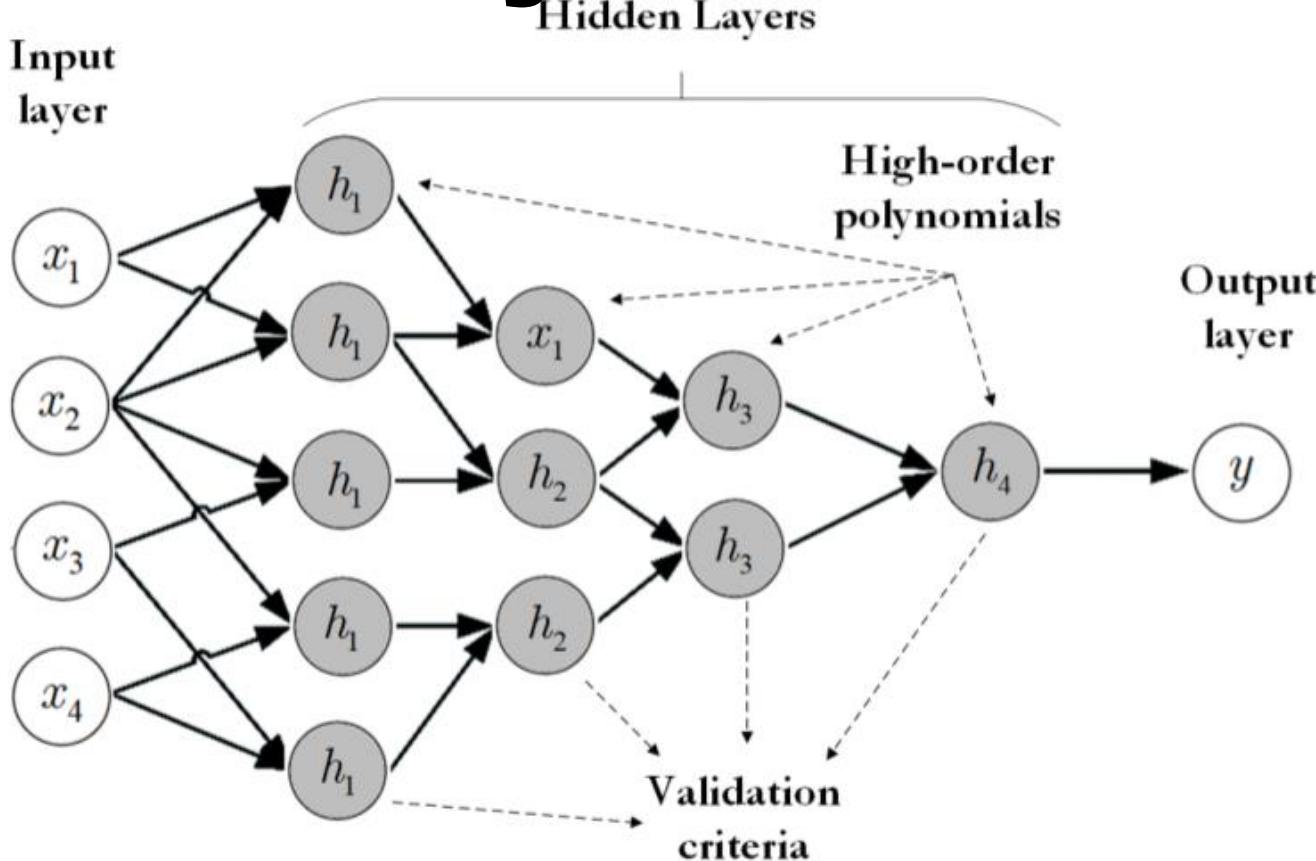
The first **artificial neural network**



Frank Rosenblatt



# A brief history... in 1966



Alexey Ivakhnenko

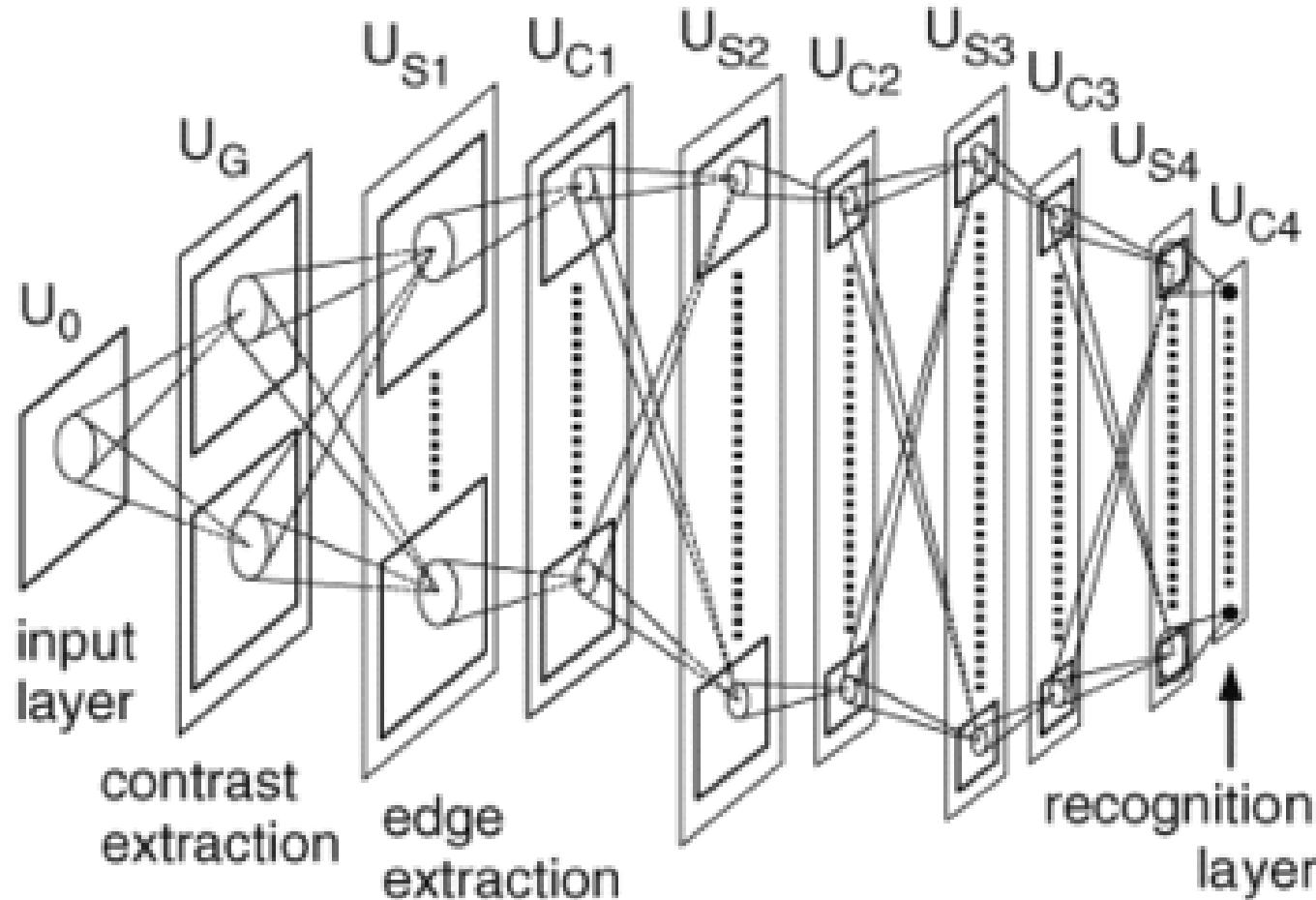
Radaideh & Kozlowski (2019)

$$y = a_0 + \sum_{i=1}^d a_i x_i + \sum_{i=1}^d \sum_{j=1}^d a_{ij} x_i x_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d a_{ijk} x_i x_j x_k + \dots$$

## Group Method of Data Handling (GMDH)

One of the first **deep (hidden-layer)** learning algorithms

# A brief history... in 1980

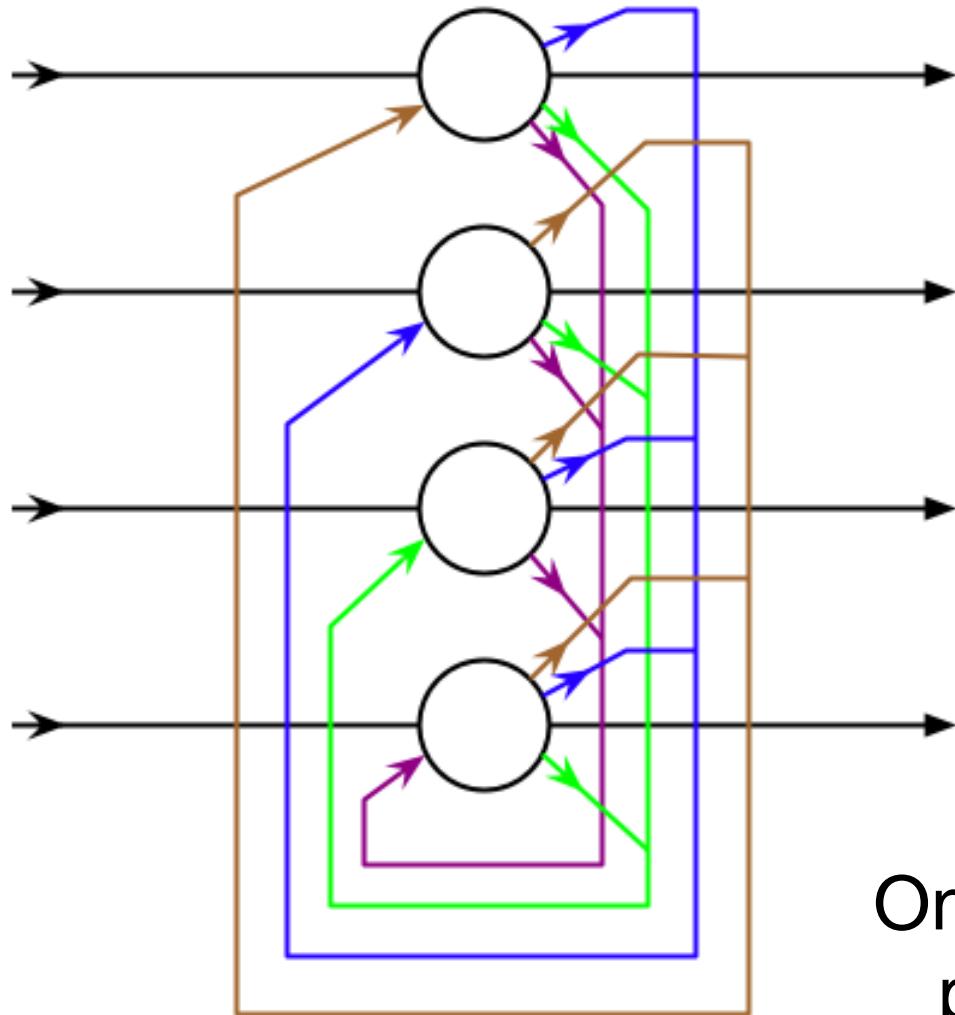


Kunihiko Fukushima

## Neocognitron

One of the first **convolutional** neural networks

# A brief history... in 1982

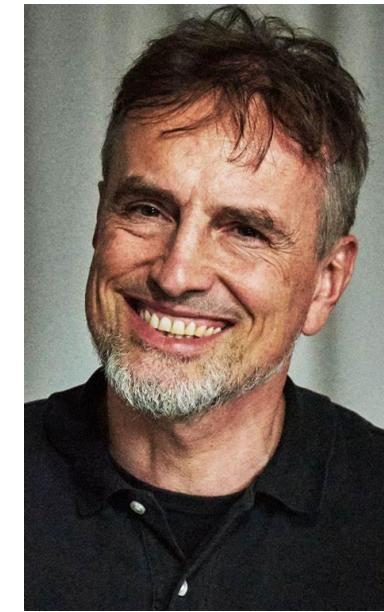
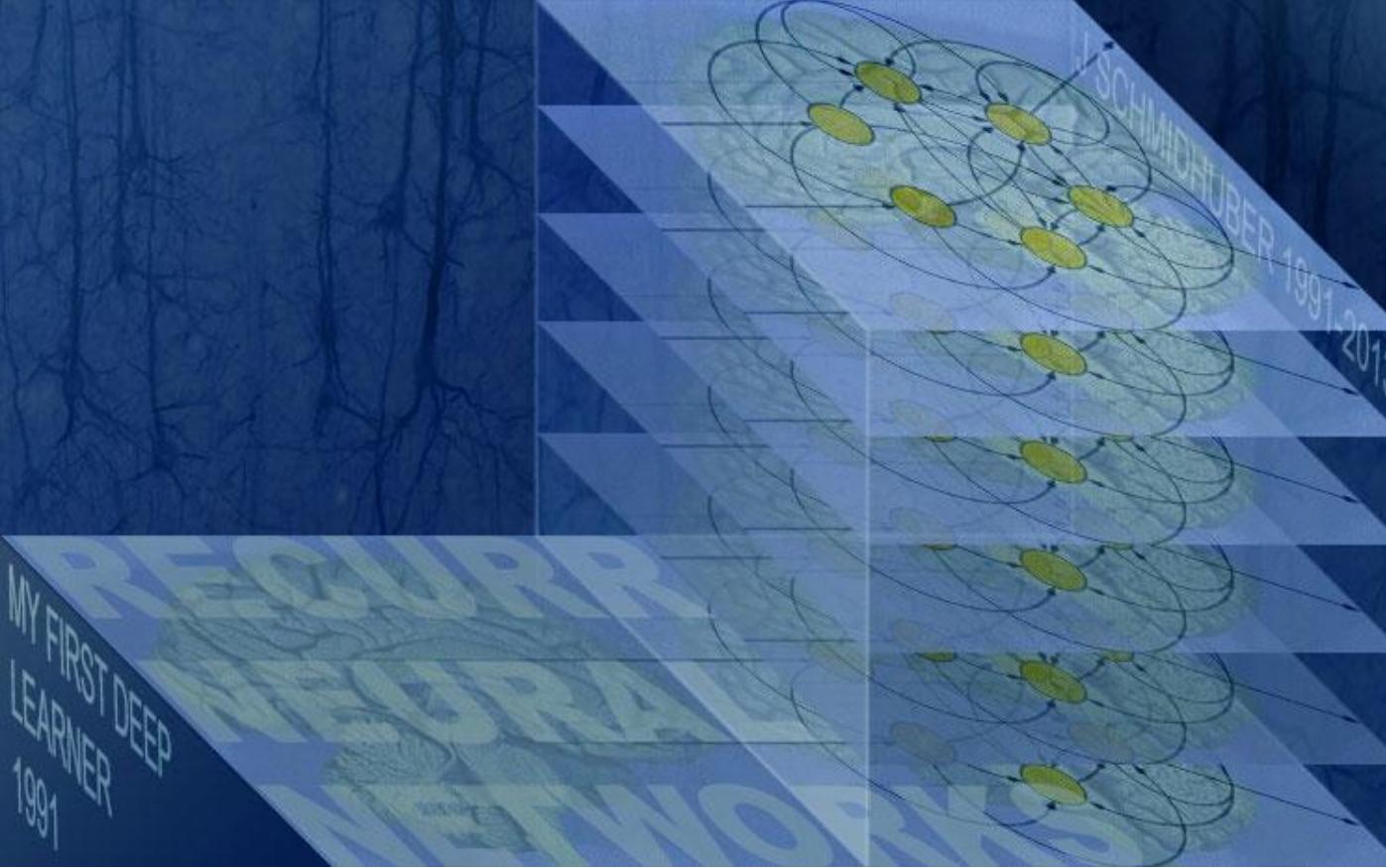


John Hopfield

## Hopfield Network

One of the first **recursive** neural networks → preserves information from former steps

# A brief history... in 1991



Jurgen Schmidhuber

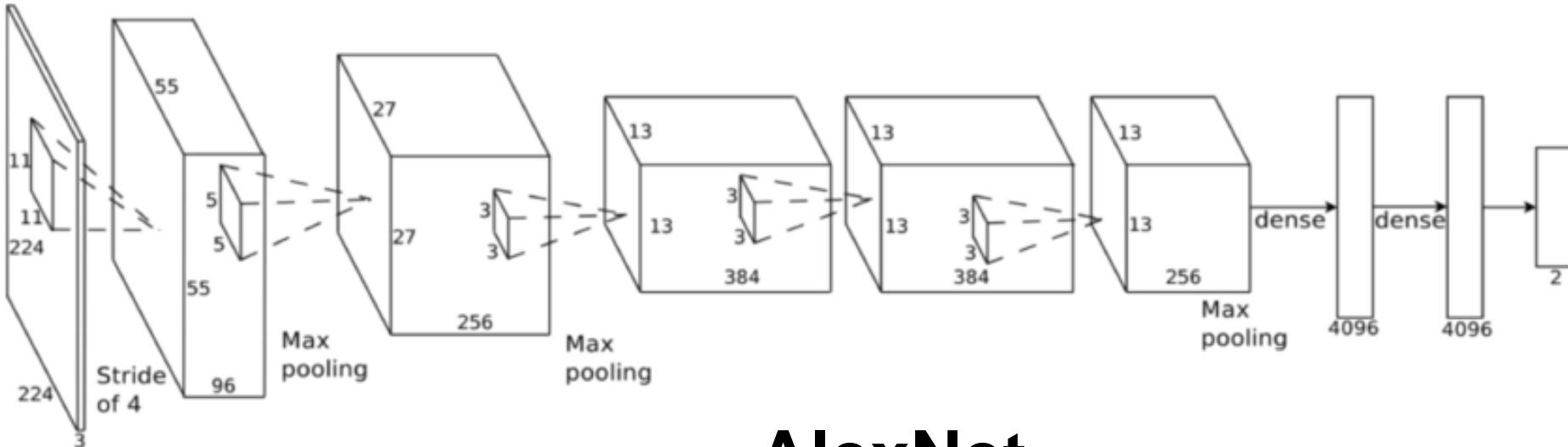
**Neural Sequence Chunker**  
Very deep learning with > 1000 layers

# A brief history... in 2011

IMAGENET



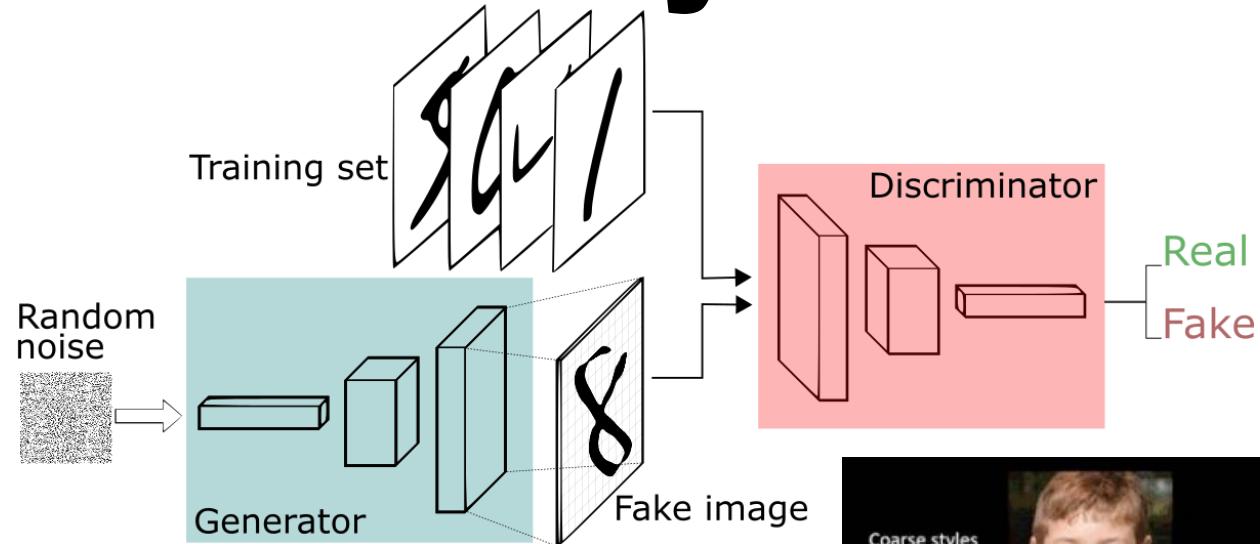
Alex Krizhevsky



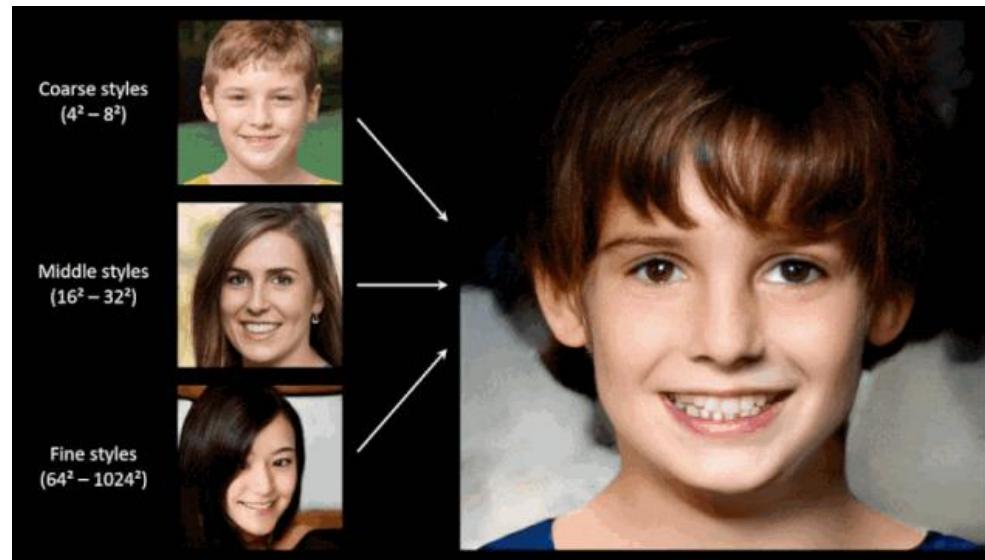
AlexNet

Putting (convolutional) **deep learning** as a superior algorithm

# A brief history... in 2014

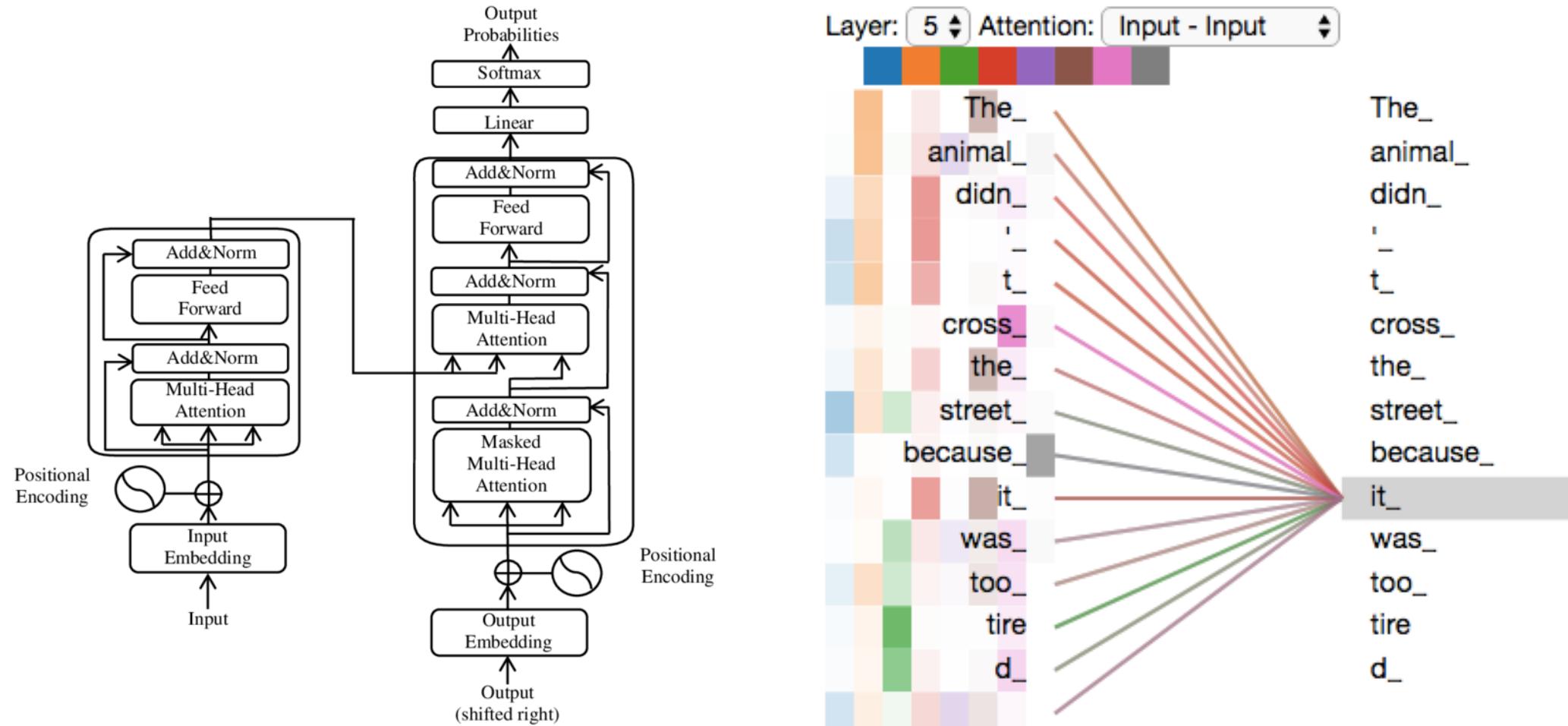


Ian Goodfellow



Generative **Adversarial** Network

# A brief history... in 2017



**Attention Is All You Need**  
Transformer model → GPT

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.



Let's do it!

## 엔드 투 엔드 오픈소스 머신러닝 플랫폼

TensorFlow 자비스크립트용 모바일 및 IoT용 프로덕션용

ML 모델을 개발하고 학습시키는 데 도움이 되는 핵심 오픈소스 라이브러리. 브라우저에서 Colab 노트북을 직접 실행하여 빠르게 시작해보세요.

TensorFlow 시작하기



# Tensorflow

(<http://www.tensorflow.org/>)

PyTorch

Get Started

Ecosystem

Mobile

Blog

Tutorials

Docs

Resources

Github

Q

# FROM RESEARCH TO PRODUCTION

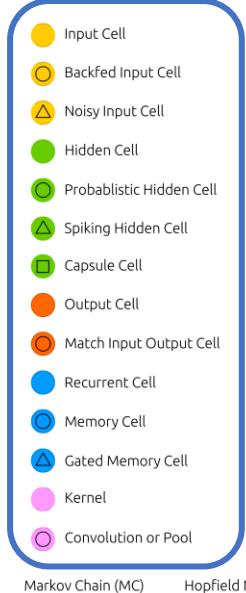
An open source machine learning framework that accelerates the path from research prototyping to production deployment.

Get Started >

Click Here to Read About Latest Updates and Improvements to PyTorch Tutorials

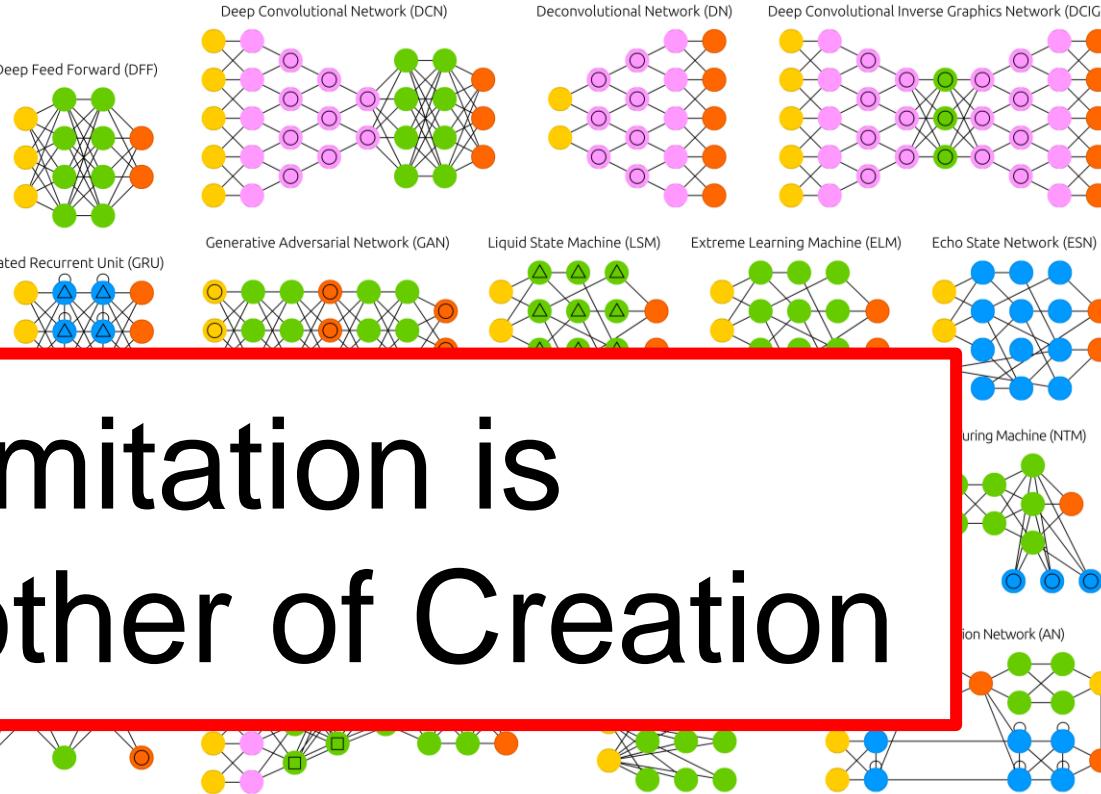
# PyTorch

(<http://pytorch.org/>)



# A mostly complete chart of Neural Networks

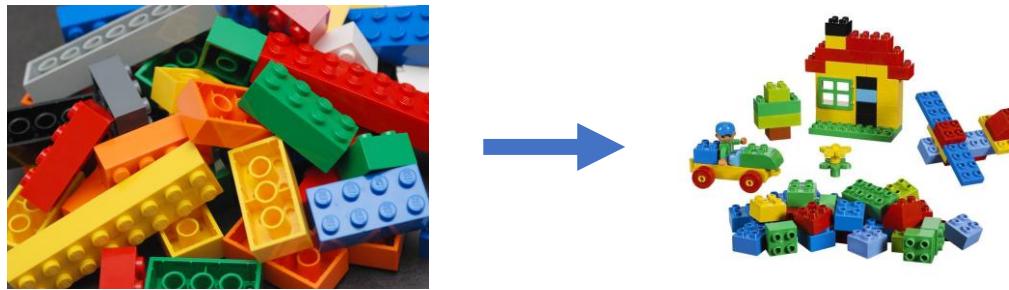
©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

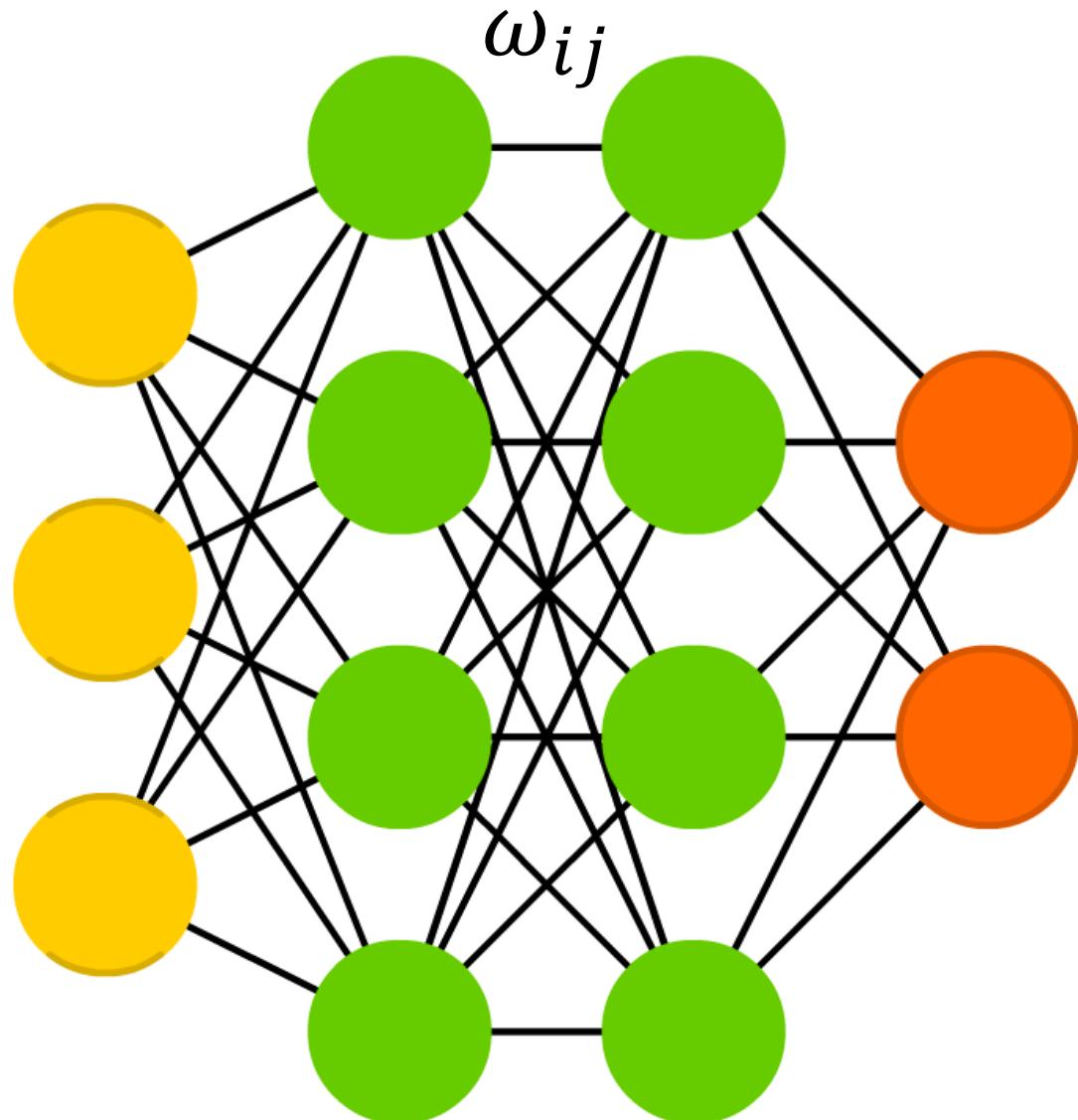


Imitation is  
the mother of Creation

<https://www.asimovinstitute.org/neural-network-zoo/>

Layer = Lego block





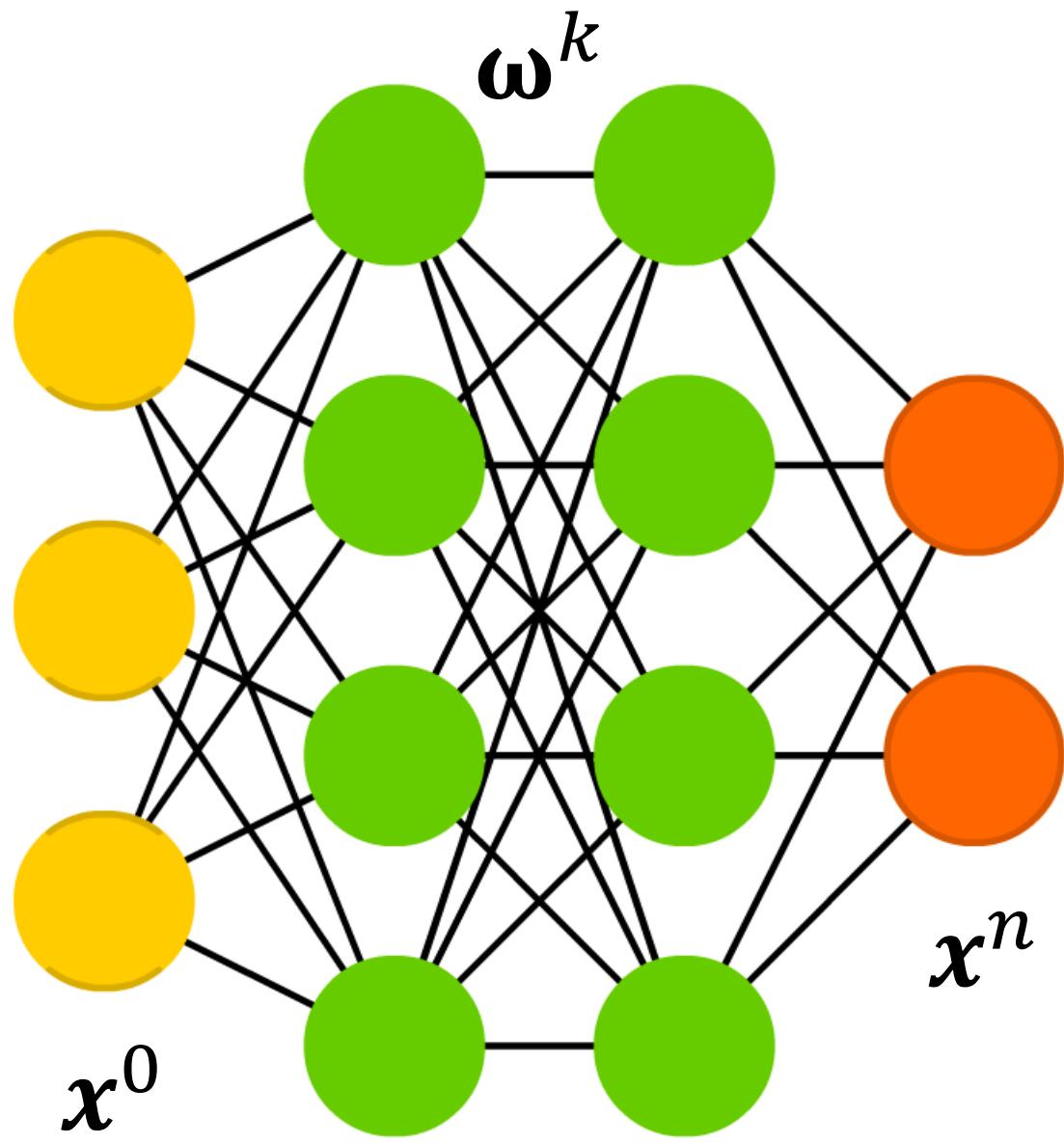
`torch.nn.Linear(n_element_before,  
n_element_after)`

## Fully Connected Layer

- Calculate the vector by using linear algebra from the vector at the former layer

$$y_i = \sum_j \omega_{ij} x_j + b_i$$

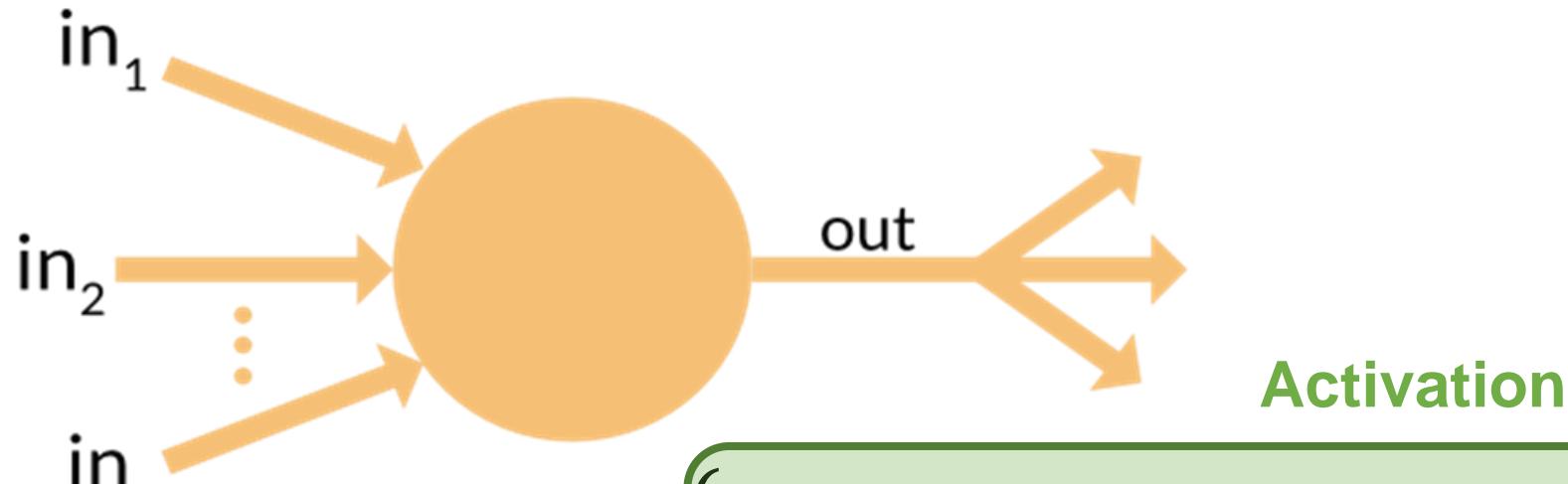
- Ideal if each vector element has a different meaning (e.g., different physical parameters)
- Useful even when different vector elements share a similar meaning (to reduce the contamination by noise)



For any set of  $\{\omega^k, b^k\}$  with  
 $x^k = \omega^k x^{k-1} + b^k$   
 $(k = 1, \dots, n),$

there **always** exists  $\{\omega, b\}$   
that gives the same result with  
 $x^n = \omega x^0 + b$

**What's the benefit of  
deep layers?**



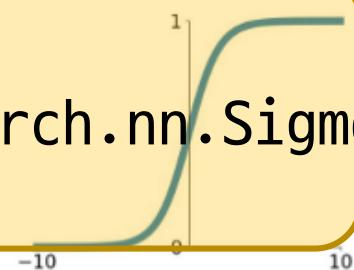
$$Y = \begin{cases} 1, & \sum_{i \in E} X_i \geq \Theta \text{ and } \sum_{j \in I} X_j = 0 \\ 0, & \sum_{i \in E} X_i < \Theta \text{ or } \sum_{j \in I} X_j > 0 \end{cases}$$

Let's apply **nonlinear** functions for each layer!  
**→ Activation Function**

# Activation Functions

## Sigmoid

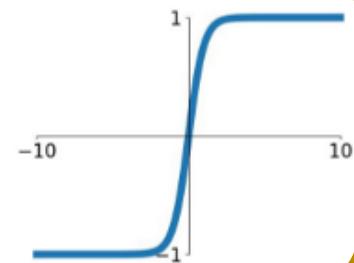
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



`torch.nn.Sigmoid()`

## tanh

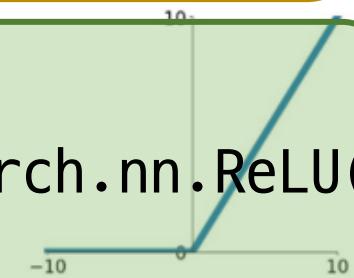
$$\tanh(x)$$



## ReLU

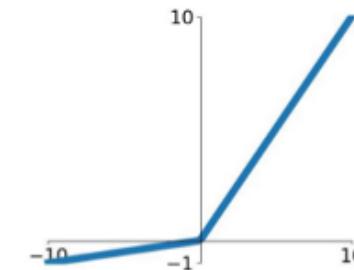
$$\max(0, x)$$

`torch.nn.ReLU()`



## Leaky ReLU

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

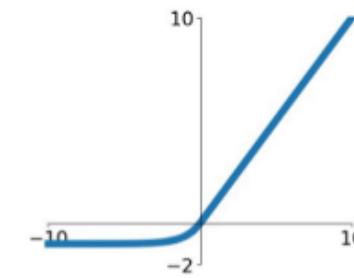


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

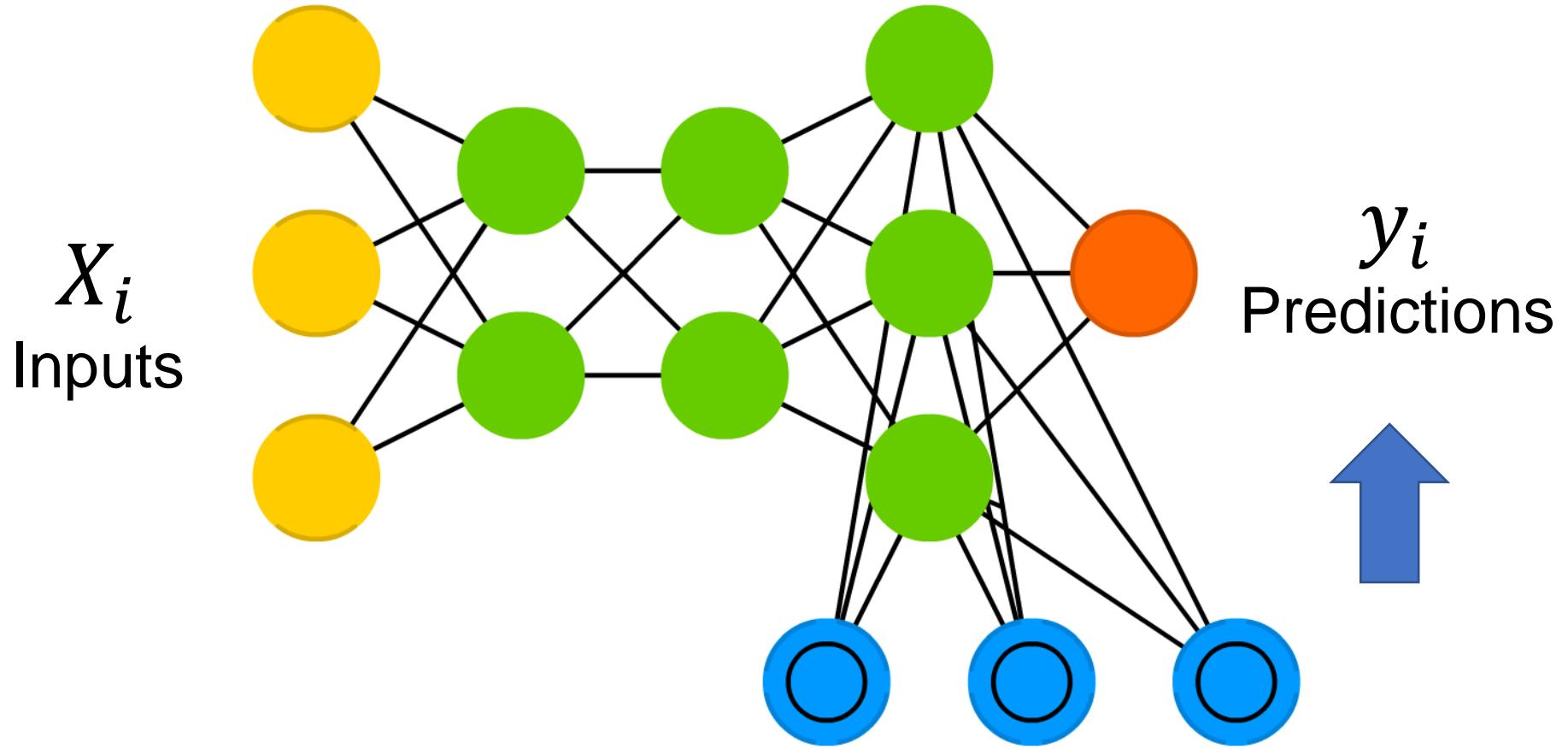
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



<https://medium.com/@shrutijadon10104776/survey-on-activation-functions-for-deep-learning-9689331ba092>

$$\textbf{Softmax } y_i = \frac{\exp x_i}{\sum_i \exp x_i} \rightarrow \sum_i y_i = 1 \quad \text{torch.nn.Softmax(dim)}$$

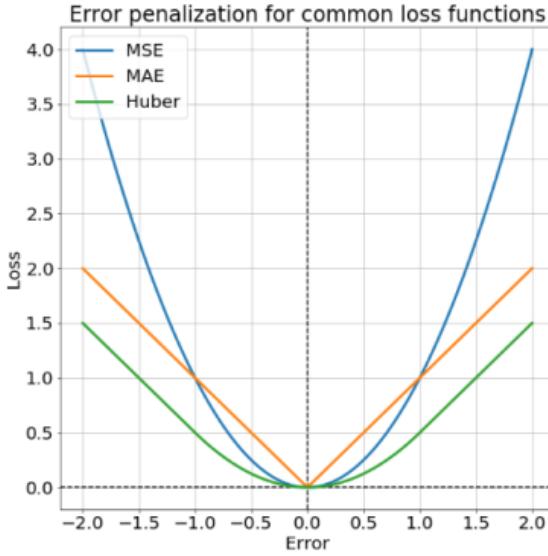
Mostly used at the last layer in classification problems



$y_i \approx Y_i$   
 $\rightarrow \text{Minimize } \mathcal{L}(\{y_i\}, \{Y_i\})$

Loss Function

# Regression



- **MSE** Mean square error (L2)
  - Strongly penalizes outliers
  - High sensitivity near minimum
- **MAE** Mean absolute error (L1)
  - Scales linearly with size of error
  - Low sensitivity near minimum
- **Huber**
  - Similar to MSE near minimum
  - Similar to MAE away from minimum

$$\mathcal{L} = \sum_i (y_i - Y_i)^2$$

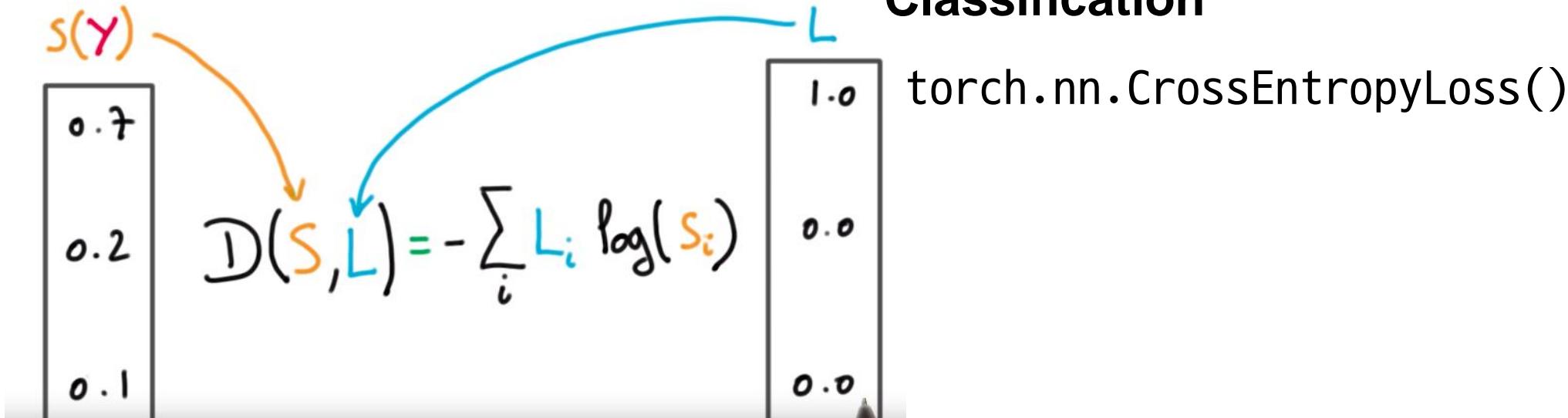
`torch.nn.MSELoss()`

$$\mathcal{L} = \sum_i |y_i - Y_i|$$

`torch.nn.L1Loss()`

<https://datascience103579984.wordpress.com/2020/01/26/introduction-to-tensorflow-in-python-from-datacamp/2/>

## CROSS-ENTROPY



## Classification

`torch.nn.CrossEntropyLoss()`

<https://medium.com/data-science-bootcamp/understand-cross-entropy-loss-in-minutes-9fb263caee9a>



# Before you start ML/DL...

- Is it really necessary?
  - Do you want to **increase precision**, or **understand the nature**?
  - Isn't classic fitting enough?
- Can you trust ML/DL results?
  - Peers/referees do not easily believe both your methods and results!
  - Please perform other analyses to convince them!



# Input/Output Data

- Which parameters will you use?
- **Normalization** (0~1? Logarithmic scale? Probability?)



**Parameters showing the strength of my character**

- Level (Lv. 1 ~ 99)
- Level / Max. Level (0 ~ 1)
- Experience Point (Exp. 1 ~  $\infty$ )
- Log (Experience Point) (0 ~  $\infty$ )
- Ranking (1st ~  $\infty$ )

# Input/Output Data

- Which parameters will you use?
- **Normalization** (0~1? Logarithmic scale? Probability?)
- Data format (Scalar? 1/2/3...-dimensional vector?)



**Parameters showing the size of a galaxy**

- Magnitude
- Magnitude – Standard Magnitude
- Stellar Mass
- Log (Stellar Mass)
- Mass Function

# Infinite Limits

After explaining to a student through various lessons and examples that

$$\lim_{x \rightarrow 8^+} \frac{1}{x - 8} = \infty$$

What if they also used other numbers?

I tried to check if he really understood that, so I gave him a different example.

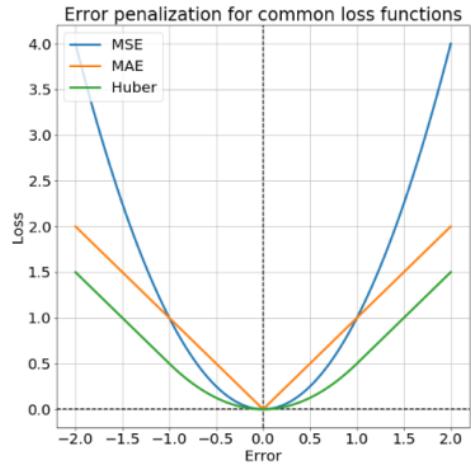
This was the result:

$$\lim_{x \rightarrow 5^+} \frac{1}{x - 5} = \text{LO}$$

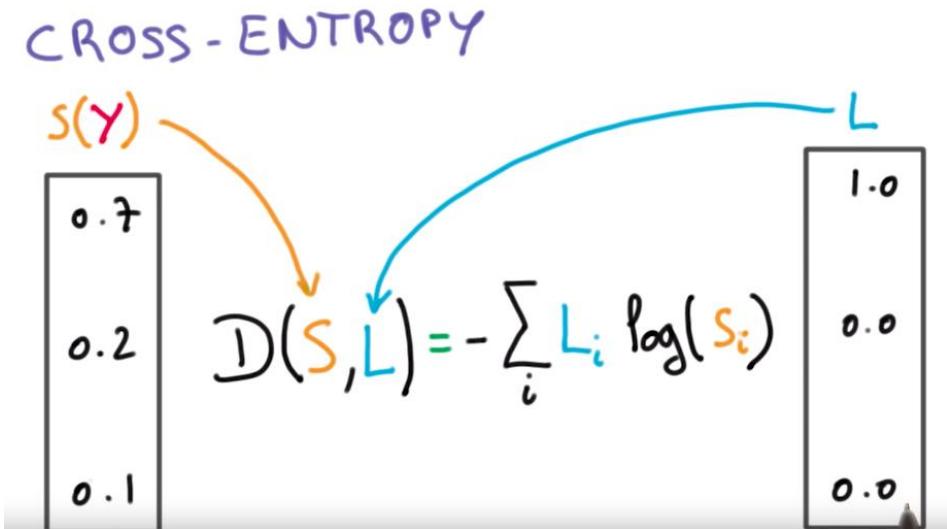
This is a math joke.....

<https://www.slideserve.com/phong/infinite-limits>

# Loss Function

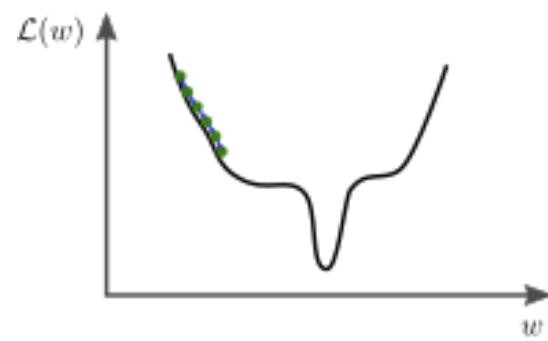


- **MSE**
  - Strongly penalizes outliers
  - High sensitivity near minimum
- **MAE**
  - Scales linearly with size of error
  - Low sensitivity near minimum
- **Huber**
  - Similar to MSE near minimum
  - Similar to MAE away from minimum

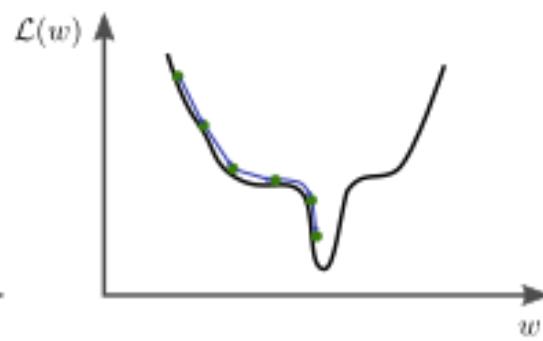


- What if there exists an unavoidable bias in your data?
- What if the given loss function fails to reproduce some important physical properties?
- What if the DL structure is too complicated?

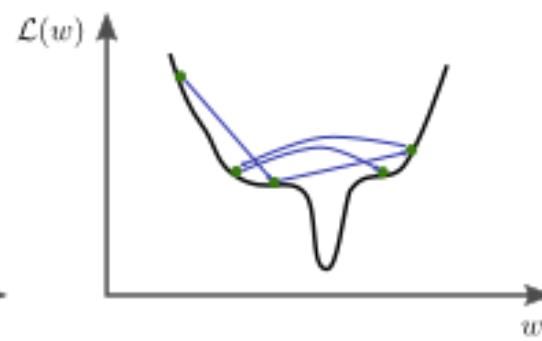
# Learning Rate



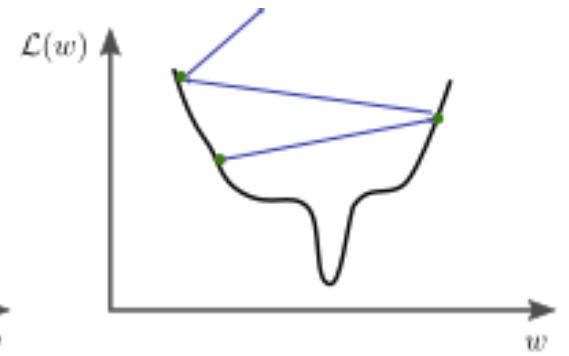
Learning rate too low



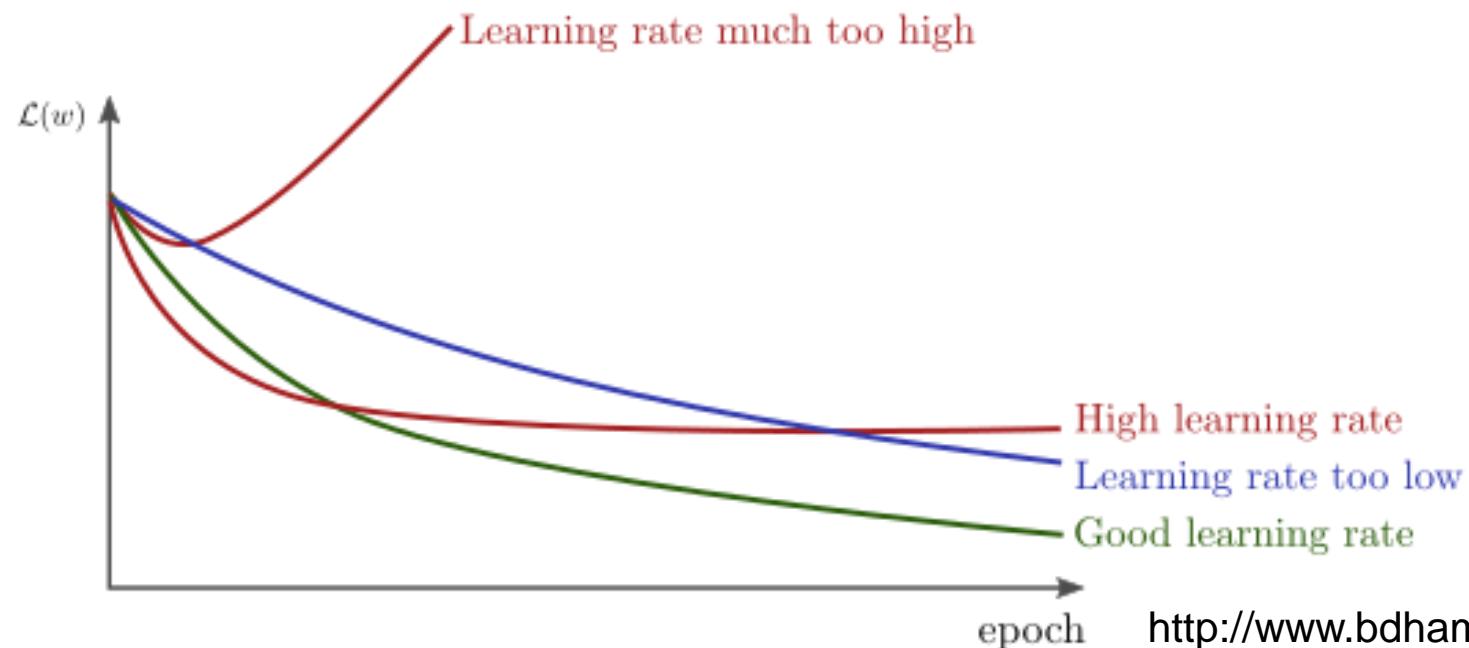
Good learning rate



High learning rate

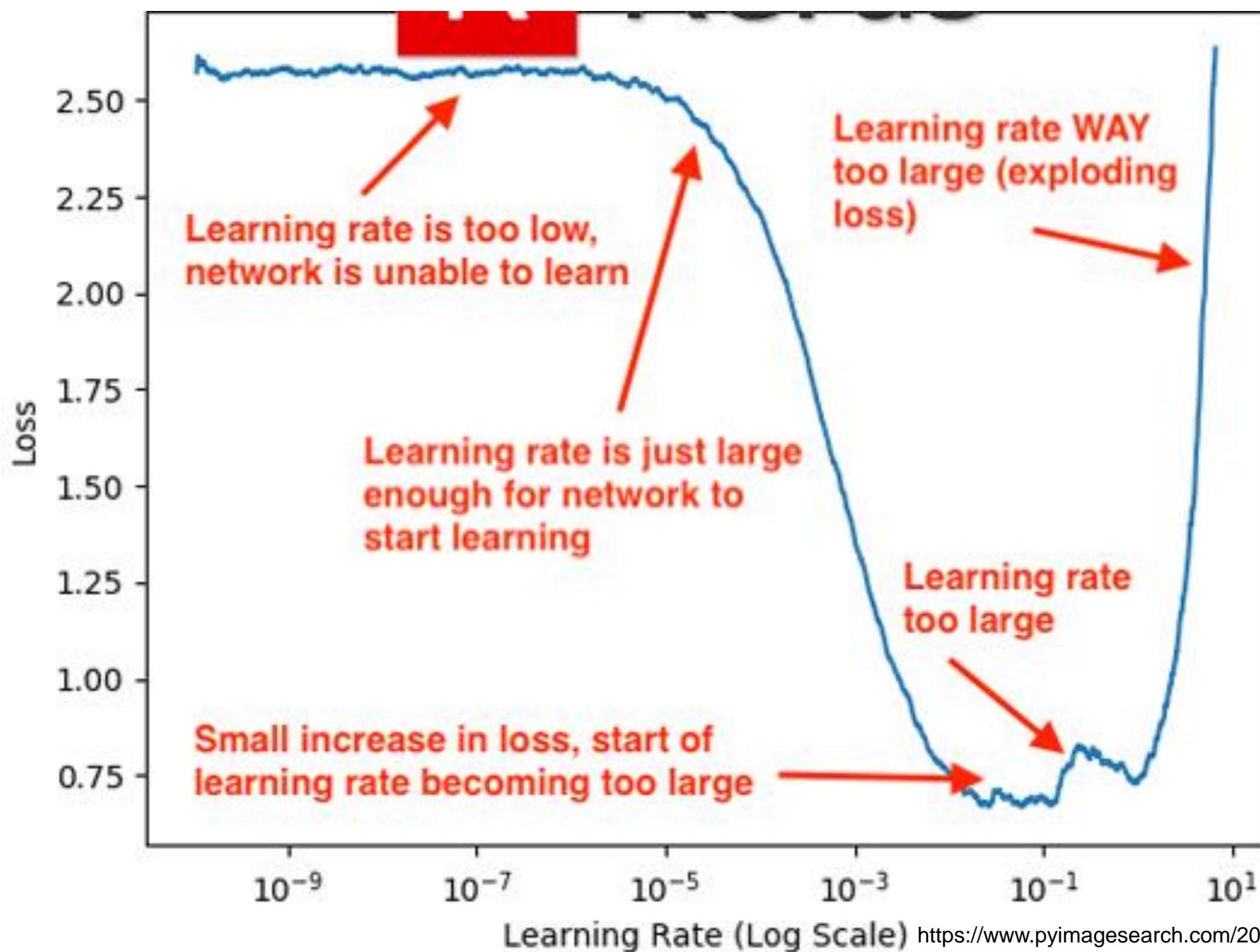


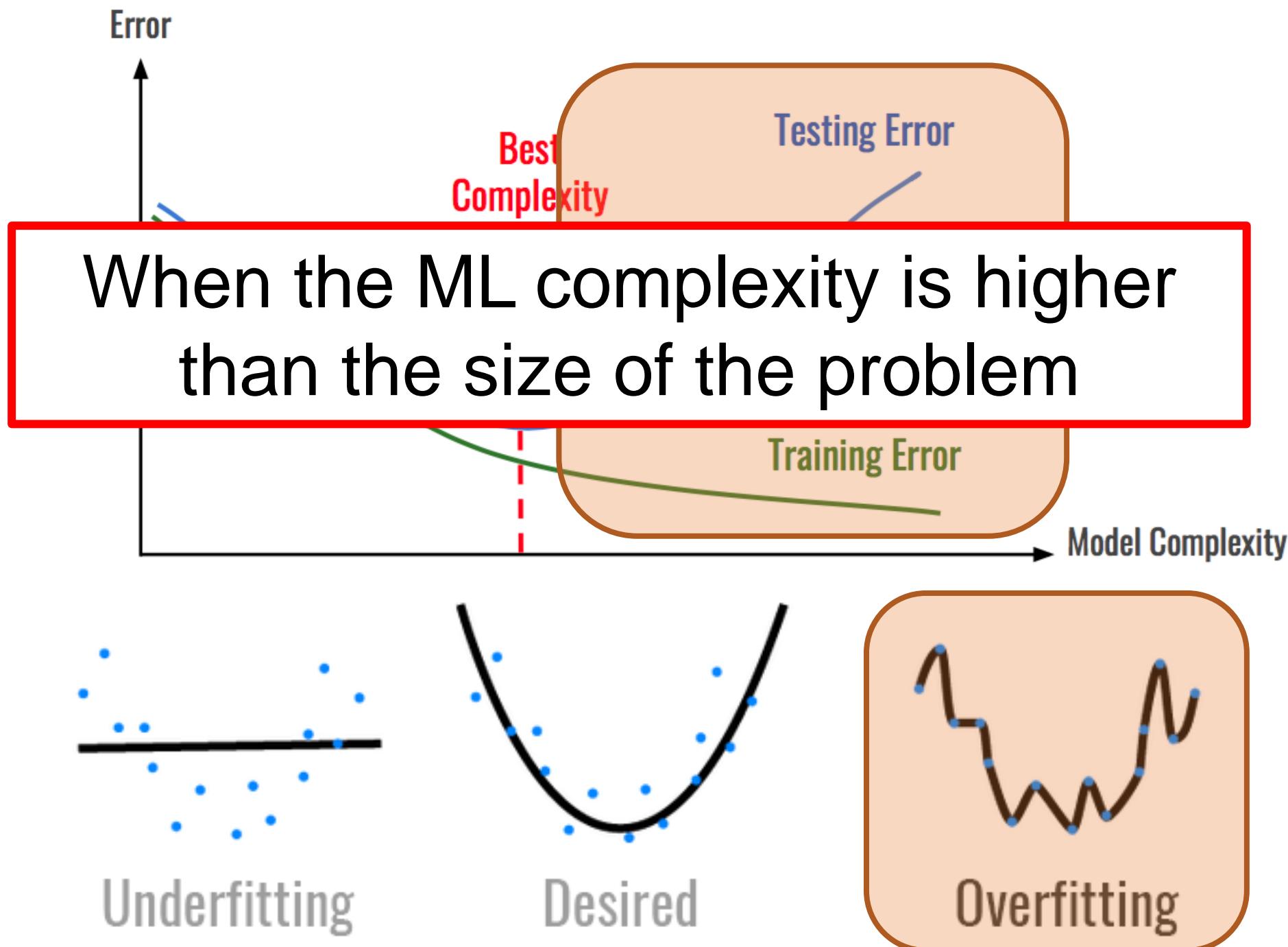
Learning rate much too high



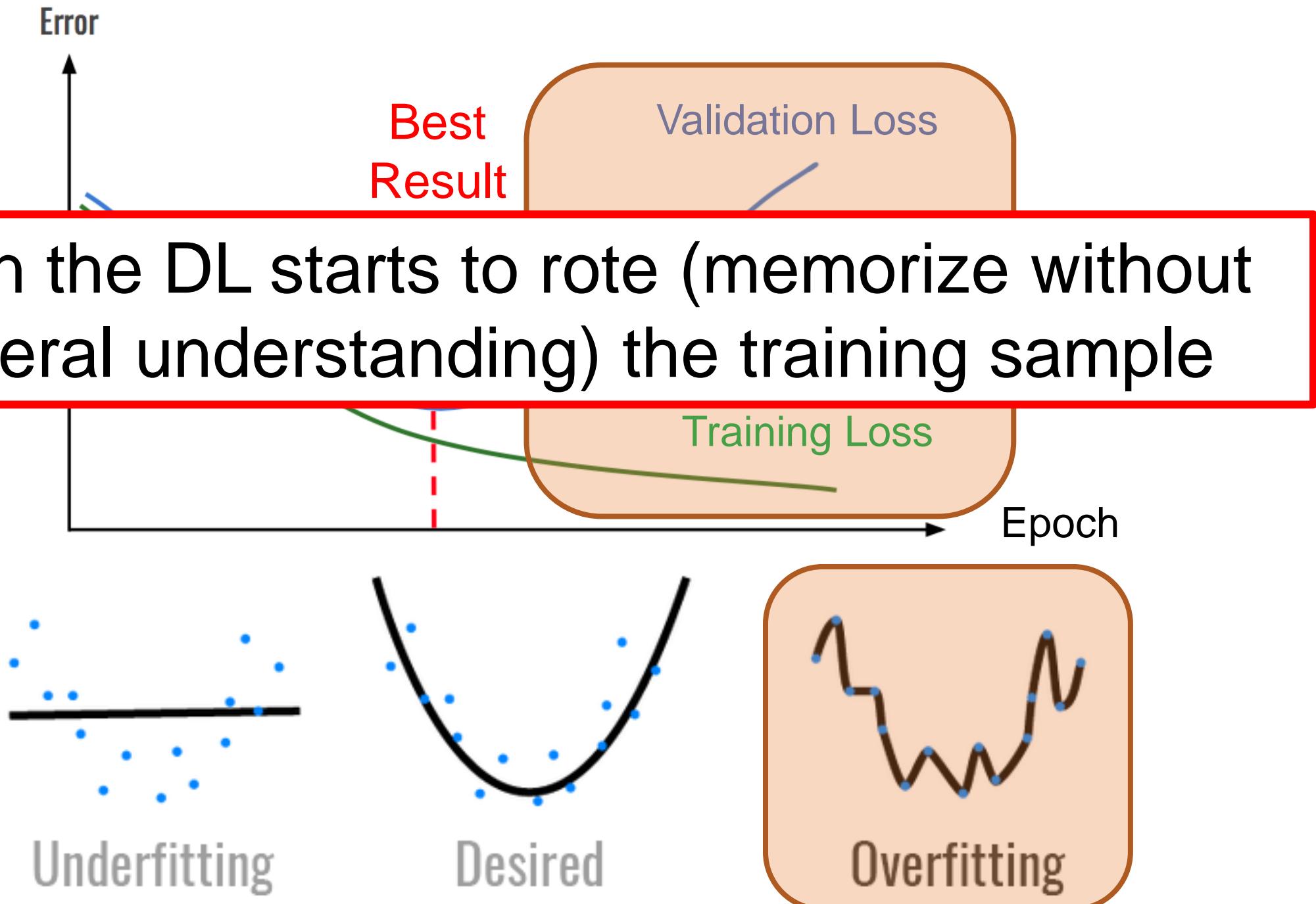
epoch

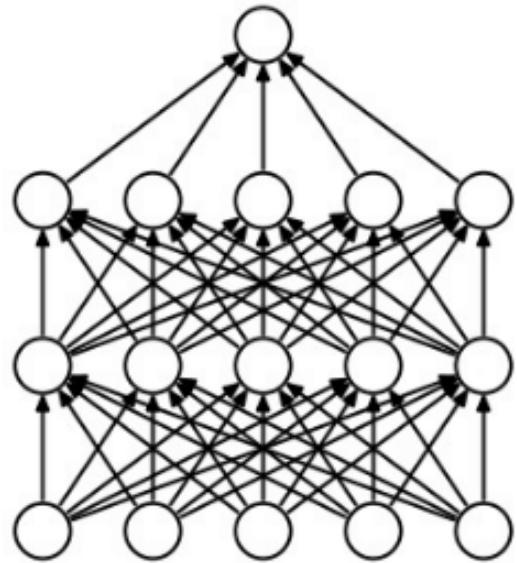
<http://www.bdhammel.com/learning-rates/>



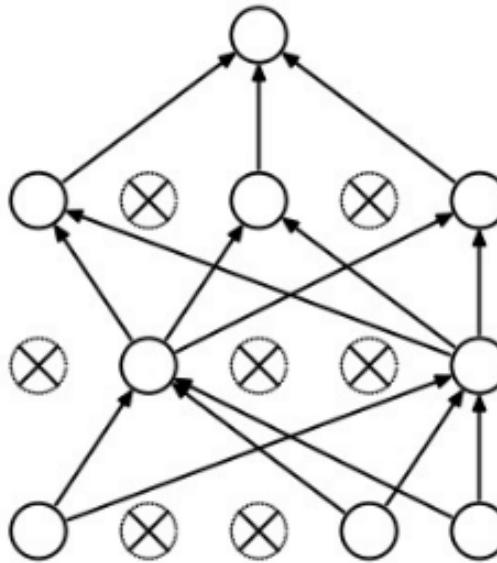


When the DL starts to rote (memorize without general understanding) the training sample





(a) Standard Neural Net



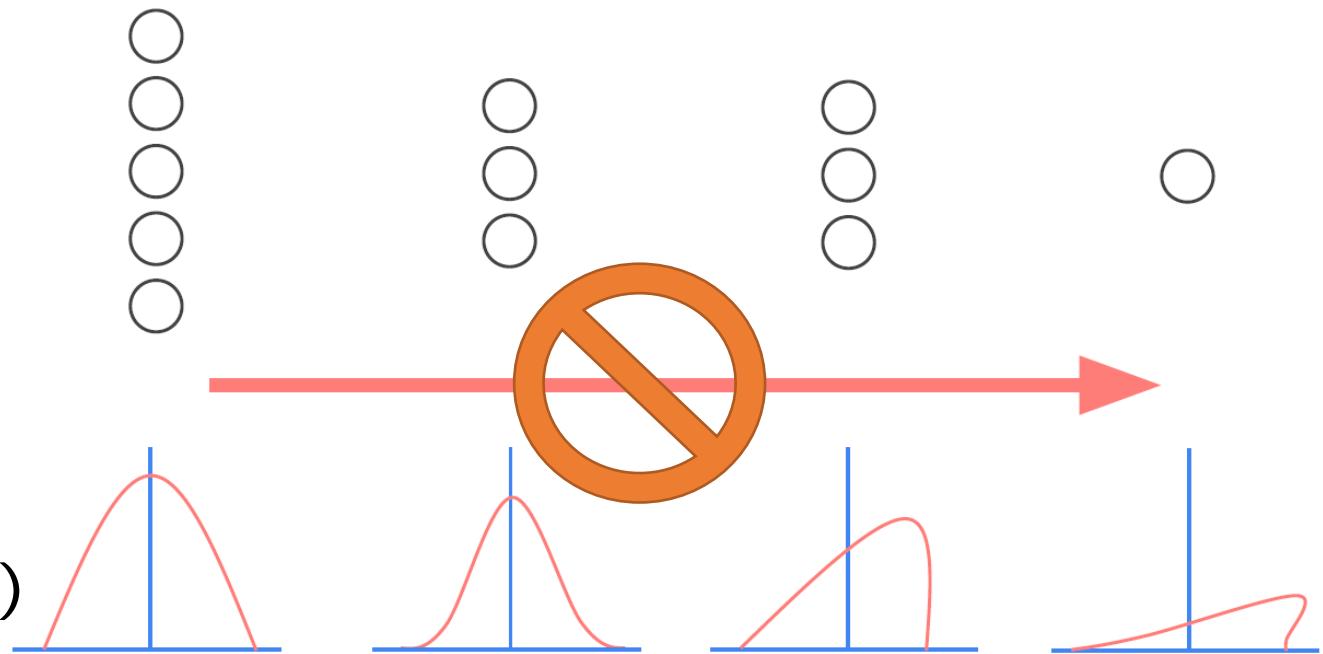
(b) After applying dropout.

# Dropout

`torch.nn.Dropout(rate)`

# Batch Normalization

`torch.nn.BatchNorm1d(num_features)`



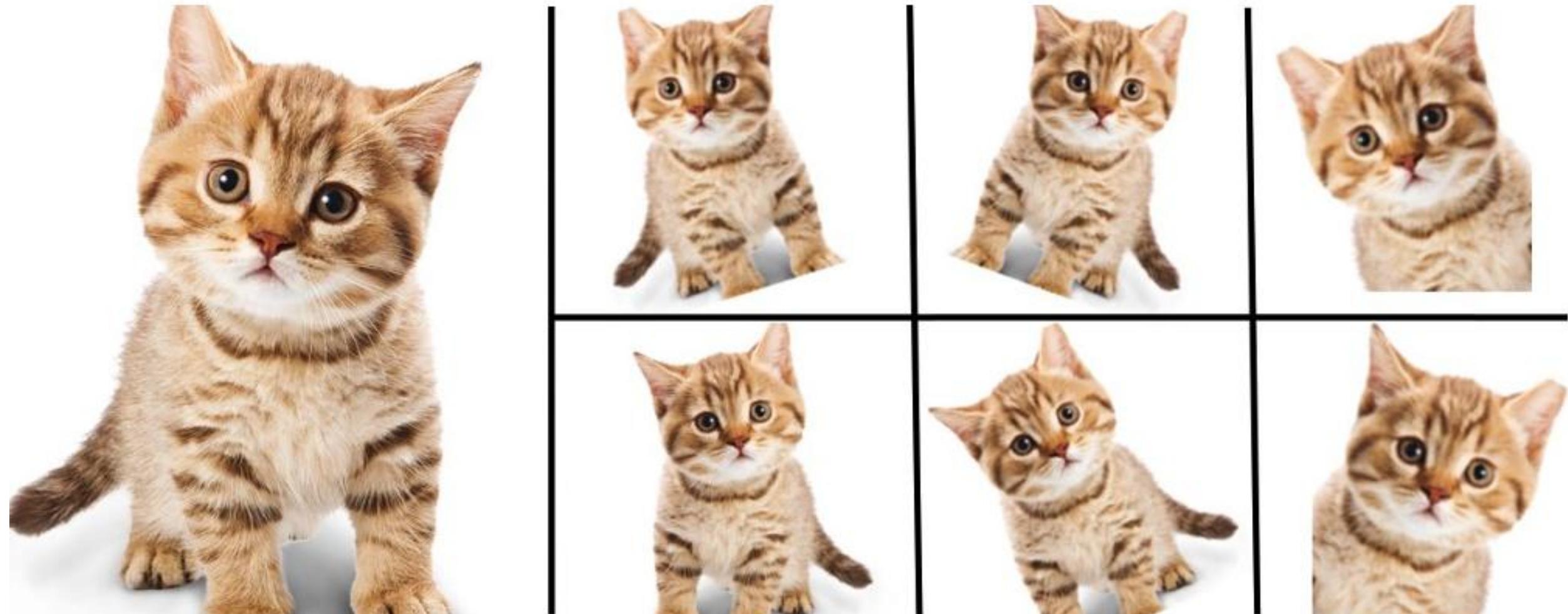


多 多 益 善

The more, the better

Paik Nam June (1987)

# Data Augmentation

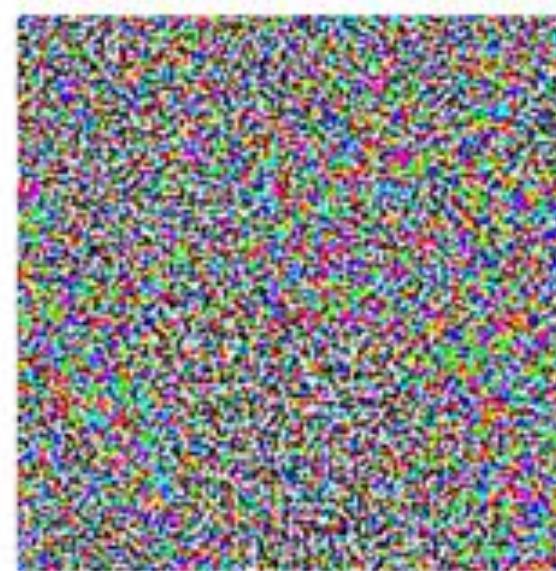


## Enlarge your Dataset



$x$   
“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

**Adding noise** may be beneficial for making ML/DL  
to concentrate only on what's important  
and forget about what's not.