

Title: Machine Learning Lecture - 230302

Speakers: Lauren Hayward

Collection: Machine Learning for Many-Body Physics (2022/2023)

Date: March 02, 2023 - 9:00 AM

URL: <https://pirsa.org/23030029>

# Machine Learning for Many-Body Physics

## Lecture 3





# State of Machine Learning and Data Science 2021

Insights from Kaggle's annual user survey focused  
on working data scientists.

**October 14, 2021**

<https://www.kaggle.com/competitions/kaggle-survey-2021/>



This is our 5th year conducting an in-depth user survey & publicly sharing the results.

Over 25,000 data scientists and ML engineers submitted responses on their backgrounds and day to day experience – everything from educational details to salaries to preferred technologies and techniques.



- **Respondents were asked about most commonly used algorithms**

- **Options included:**

- ▶ Linear or logistic regression
- ▶ Decision trees or random forests
- ▶ Gradient boosting machines
- ▶ Convolutional neural networks
- ▶ Bayesian approaches
- ▶ Dense neural networks
- ▶ Recurrent neural networks
- ▶ Transformer networks
- ▶ Generative adversarial networks
- ▶ Evolutionary approaches

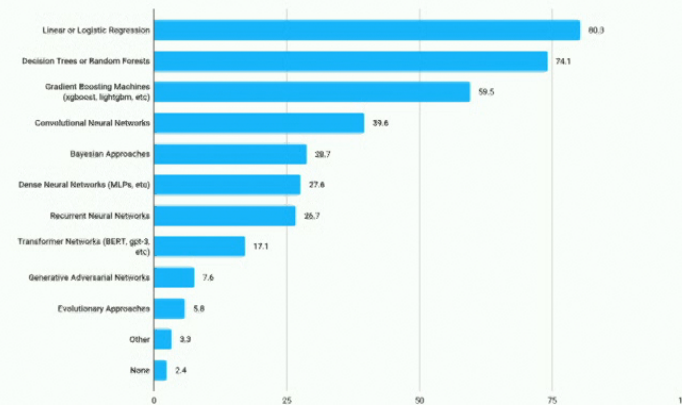
<https://www.kaggle.com/competitions/kaggle-survey-2021/>

## Methods & Algorithms

Like last year, the most commonly used algorithms were linear and logistic regression, followed closely by decision trees and random forests.

Of more complex methods, gradient boosting machines and convolutional neural networks were the most popular approaches.

Methods and Algorithms Usage



<https://www.kaggle.com/competitions/kaggle-survey-2021/>

Recall the goal of supervised learning (SL): Given a dataset  $\mathcal{D} = \{(\vec{x}, \vec{y})\}$ , fit a function  $\vec{f}(\vec{x})$  to  $\vec{y}$

$$\vec{x} = (x_1, x_2, \dots, x_{d_x})$$

$$\vec{y} = (y_1, y_2, \dots, y_{d_y})$$

$$\vec{f} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$$

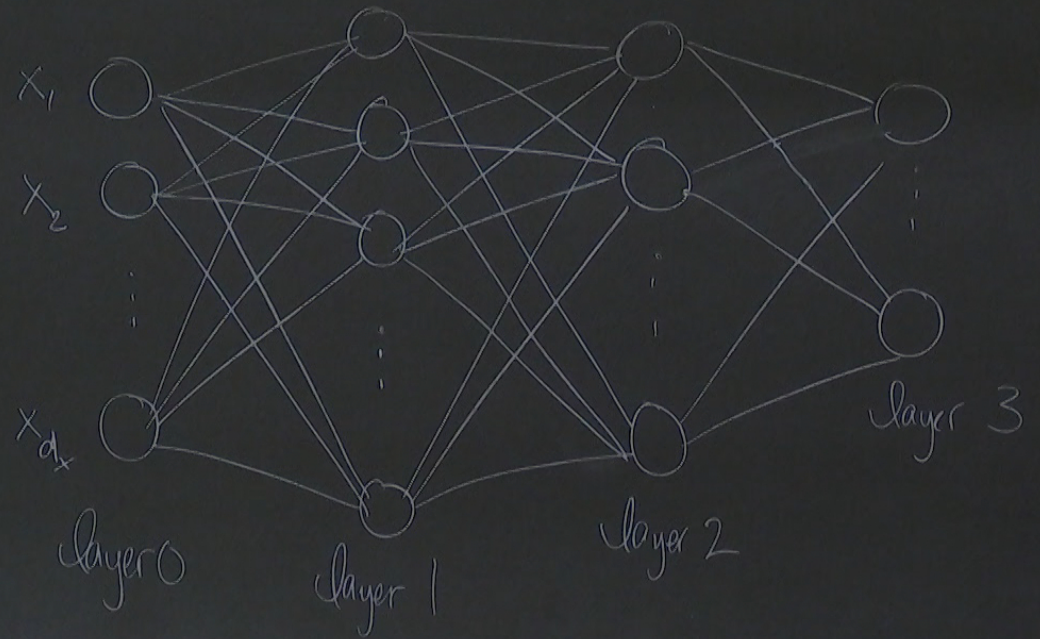
$$: \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$$

Outline for today: intro to SL with feedforward neural networks (NNs)

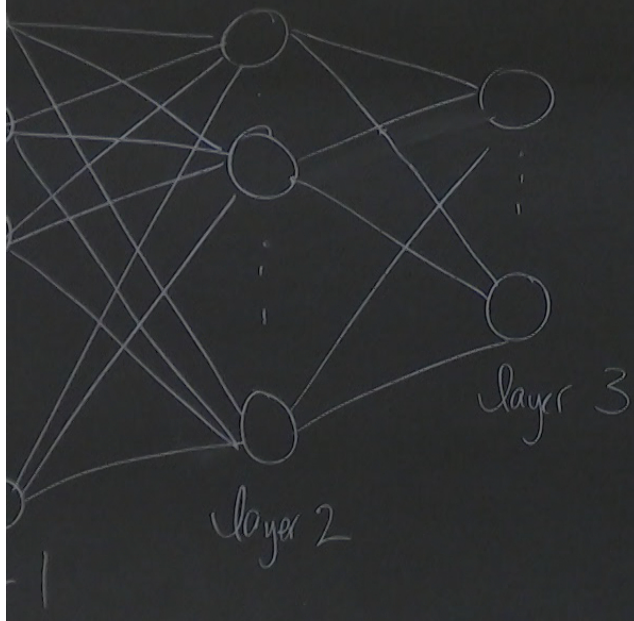
- Network architecture and propagation of information
  - ↳ weights
  - ↳ biases
  - ↳ activation functions
- Expressivity of NNs
- Basics of training ("learning")
  - ↳ cost functions
  - ↳ learning algorithms



# Feedforward NNs



• Basics of training ("learn")  
↳ cost functions  
↳ learning algorithms



Each  $\bigcirc$  is a "neuron"

Notation.

- $n_l \equiv \#$  of neurons in layer  $l$
- $L \equiv$  largest value of  $l$  (eg.  $L=3$  here)
- $a_j^{(l)} \equiv$  output from the  $j^{\text{th}}$  neuron in layer  $l$

Then:

$$n_0 = d_x$$

$$n_L = d_y$$

$$a_i^{(0)} = x_i \quad (1 \leq i \leq d_x)$$

ayer 1

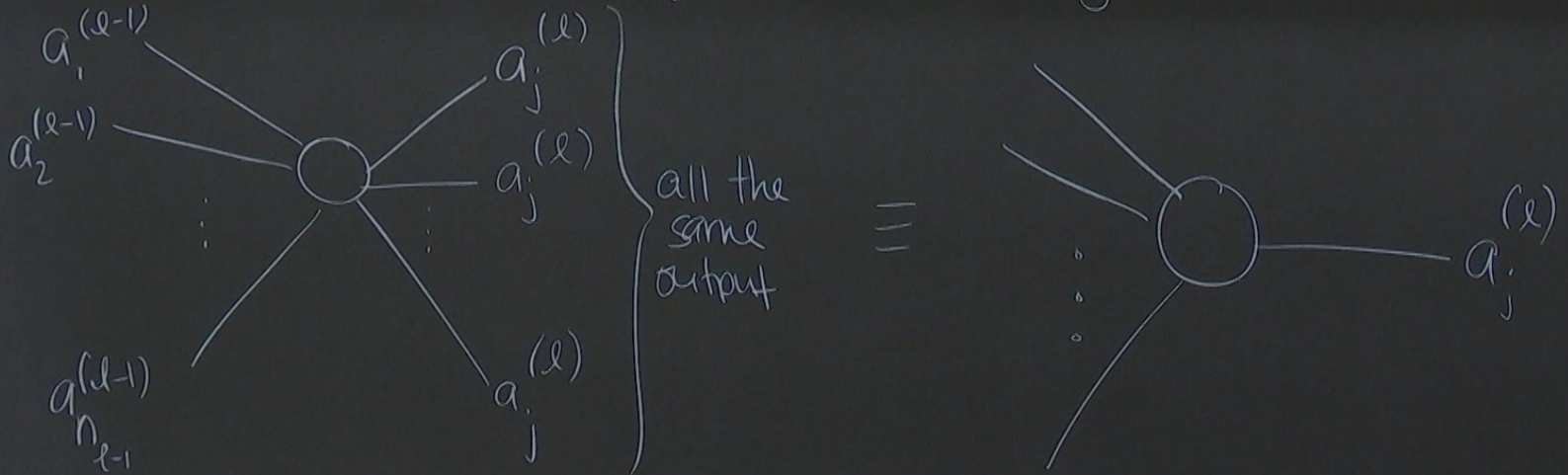
layer 0

layer 1

layer 2

$$a_j^{(l)} = \sigma(\dots)$$

Let's zoom in on the  $j^{\text{th}}$  neuron in layer  $l > 0$ :



$\equiv$  output from the  $j^{\text{th}}$  neuron in layer  $l$

The output  $a_j^{(l)}$  is.

$$a_j^{(l)} = g_l \left( \underbrace{\sum_{i=1}^{n_{l-1}} a_i^{(l-1)} W_{ij}^{(l)} + b_j^{(l)}}_{\equiv} \right)$$

The output  $a_j^{(l)}$  is:

$$a_j^{(l)} = g_l \left( \sum_{i=1}^{n_{l-1}} a_i^{(l-1)} W_{ij}^{(l)} + b_j^{(l)} \right)$$

weight  $W_{ij}^{(l)}$  for each "link"  
 bias  $b_j^{(l)}$  for each neuron  
 activation function  $g_l$   
 $\equiv z_j^{(l)}$

The weights and biases are changed as the network "learns"

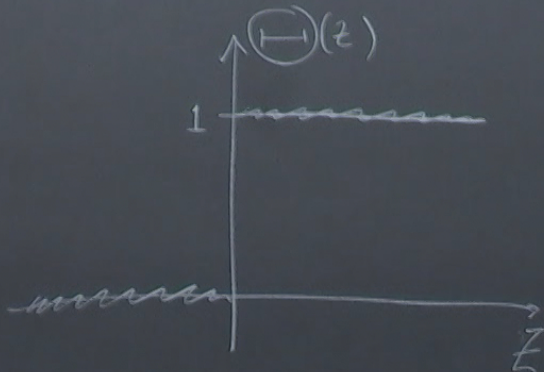
We can write the output from layer  $l$  as:

$$\vec{a}^{(l)} \equiv (a_1^{(l)}, a_2^{(l)}, \dots, a_{n_l}^{(l)})$$

Then the output from the last layer is  $\vec{a}^{(L)} = \vec{f}(\vec{x})$

Activation functions: various choices are possible

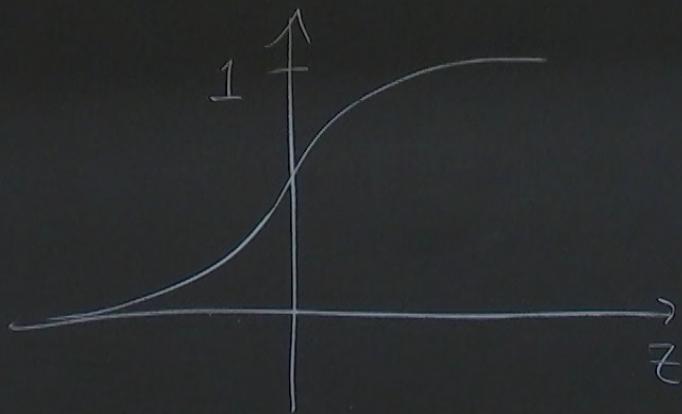
① Perceptron  $\Theta(z)$   
(step function)



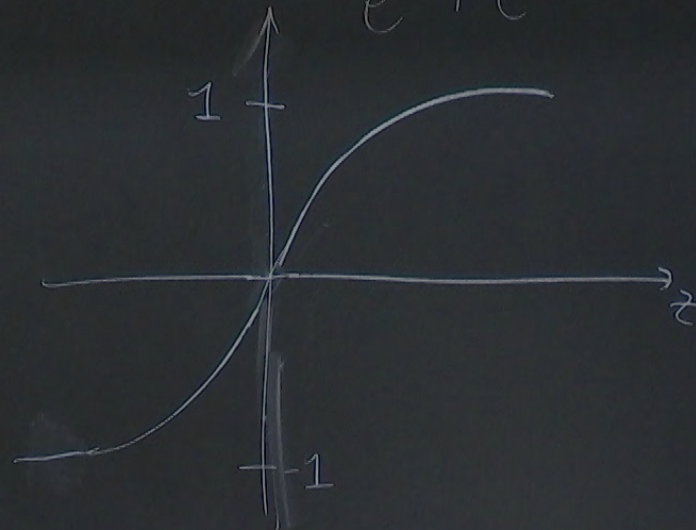
learning is difficult  
because most algorithms  
calculate the gradient  
of the neuron's output



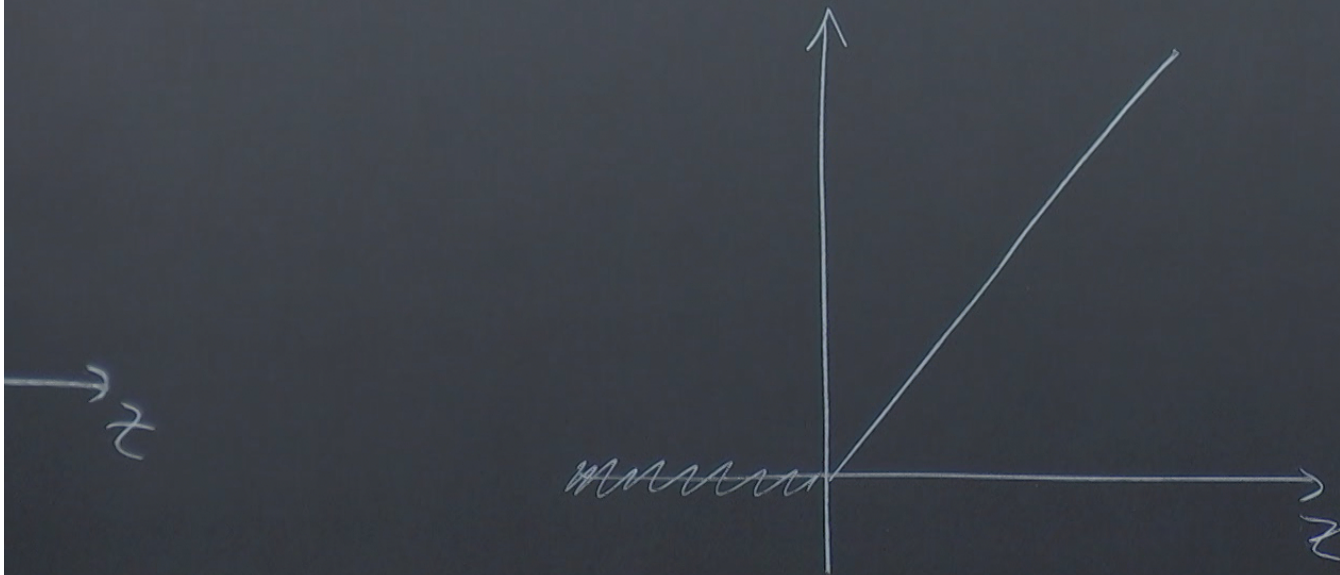
② Sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$



③  $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$



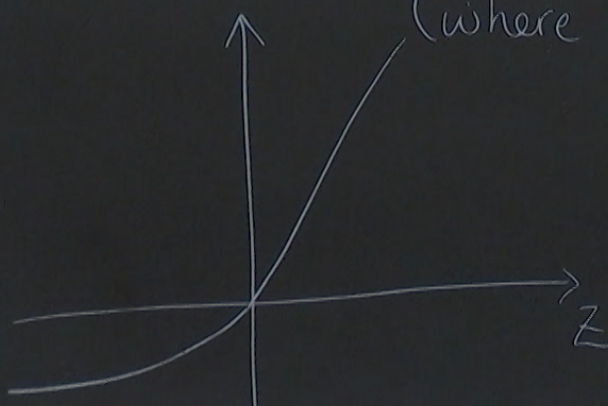
④ Rectified linear unit  
 $\text{ReLU}(z) = \max(0, z)$



⑤ Exponential linear unit

$$\text{ELU}(z) = \begin{cases} \alpha(e^z - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

(where  $\alpha > 0$ )



+1

o o o

For classification problems with  $d_y > 1$ , it is common to apply the softmax activation function to the last layer:

$$\left[ \text{softmax}(z_1, z_2, \dots, z_{n_e}) \right]_i = \frac{e^{z_i}}{\sum_{j=1}^{n_e} e^{z_j}} \quad (1 \leq i \leq n_e)$$

can interpret outputs as probabilities

$$n = d_u$$

## Expressivity of NNs (see Neilsen, Chap. 4)

Universal approximation theorem: given a continuous  $f(x)$ , there exists a NN with one hidden layer that can approximate  $f(x)$  with arbitrary accuracy

→ learning algorithms

Caveats       $\vec{x} = (x_1, x_2, \dots, x_{d_x})$

- doesn't guarantee physical interpretability
- the NN required might be huge
- knowing that such a NN exists does not guarantee we can find it
- a more expressible function can be more prone to overfitting

Caveats      $\vec{X} = (x_1, x_2, \dots, x_{d_x})$

- doesn't guarantee physical interpretability
- the NN required might be huge
- knowing that such a NN exists does not guarantee we can find it
- a more expressible function can be more prone to overfitting

→ learning algorithms

Now let's get our NN to learn!

## Cost function

The goal of SL with NNs is to find weights and biases such that for input  $\vec{x}$ , the NN output  $\vec{a}^{(L)}(\vec{x})$  is "close" to the label  $\vec{y}(\vec{x})$ .



We use a cost function to measure how well the NN approximates the labels.

Different options:

① Mean-squared error (MSE)

$$C_{\text{MSE}} = \frac{1}{2|\mathcal{D}|} \sum_{\vec{x} \in \mathcal{D}} \sum_{i=1}^{n_L} \left[ a_i^{(L)}(\vec{x}) - y_i(\vec{x}) \right]^2$$

$\underbrace{\hspace{10em}}_{\|\vec{a}^{(L)}(\vec{x}) - \vec{y}(\vec{x})\|^2}$

② Cross

② Cross entropy (CE) (for  $0 \leq y_i(\vec{x}) \leq 1$   
 $0 \leq a_i^{(L)}(\vec{x}) \leq 1$ )

$$C_{CE} = \frac{1}{|\mathcal{D}|} \sum_{\vec{x} \in \mathcal{D}} \sum_{i=1}^{n_L} \left[ y_i(\vec{x}) \log a_i^{(L)}(\vec{x}) + (1 - y_i(\vec{x})) \log (1 - a_i^{(L)}(\vec{x})) \right]$$

Next time we will introduce algorithms that aim to minimize  $C$