

# Big data science Day 4

F. Legger - INFN Torino

<https://github.com/Course-bigDataAndML/MLCourse-2223>

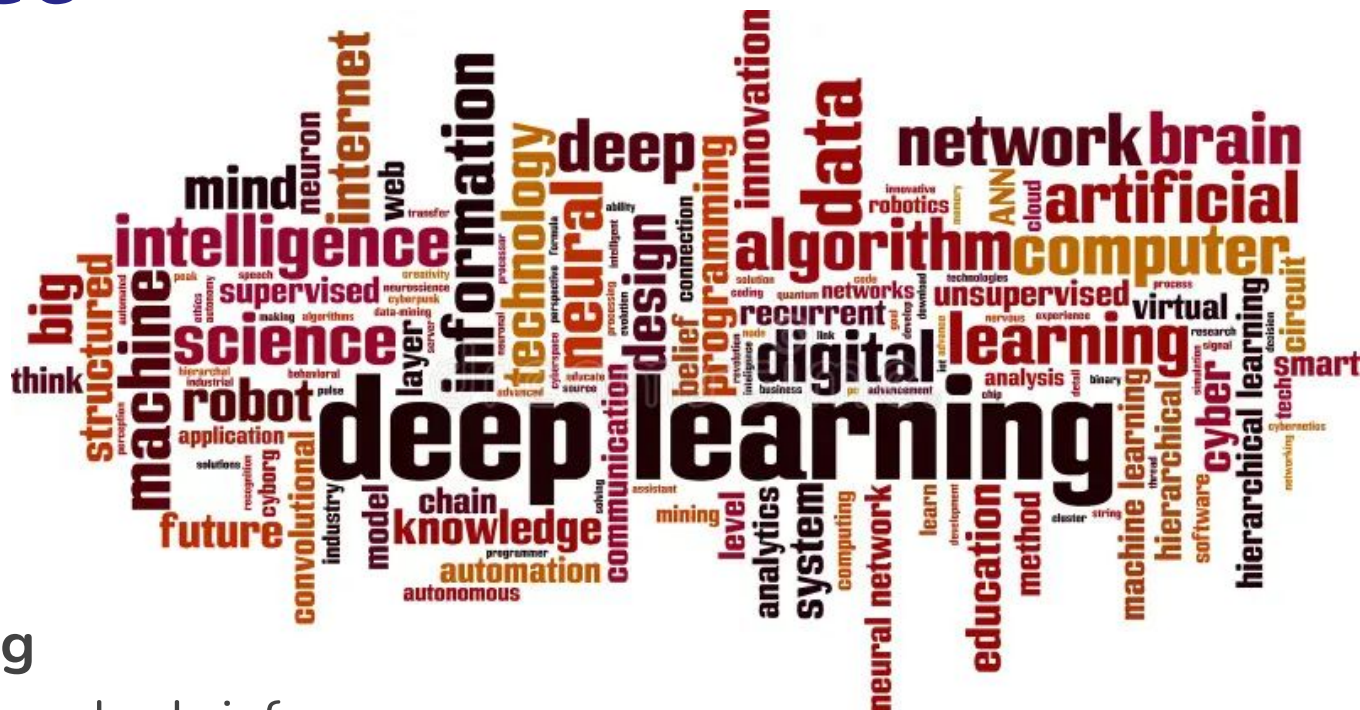


# We learned

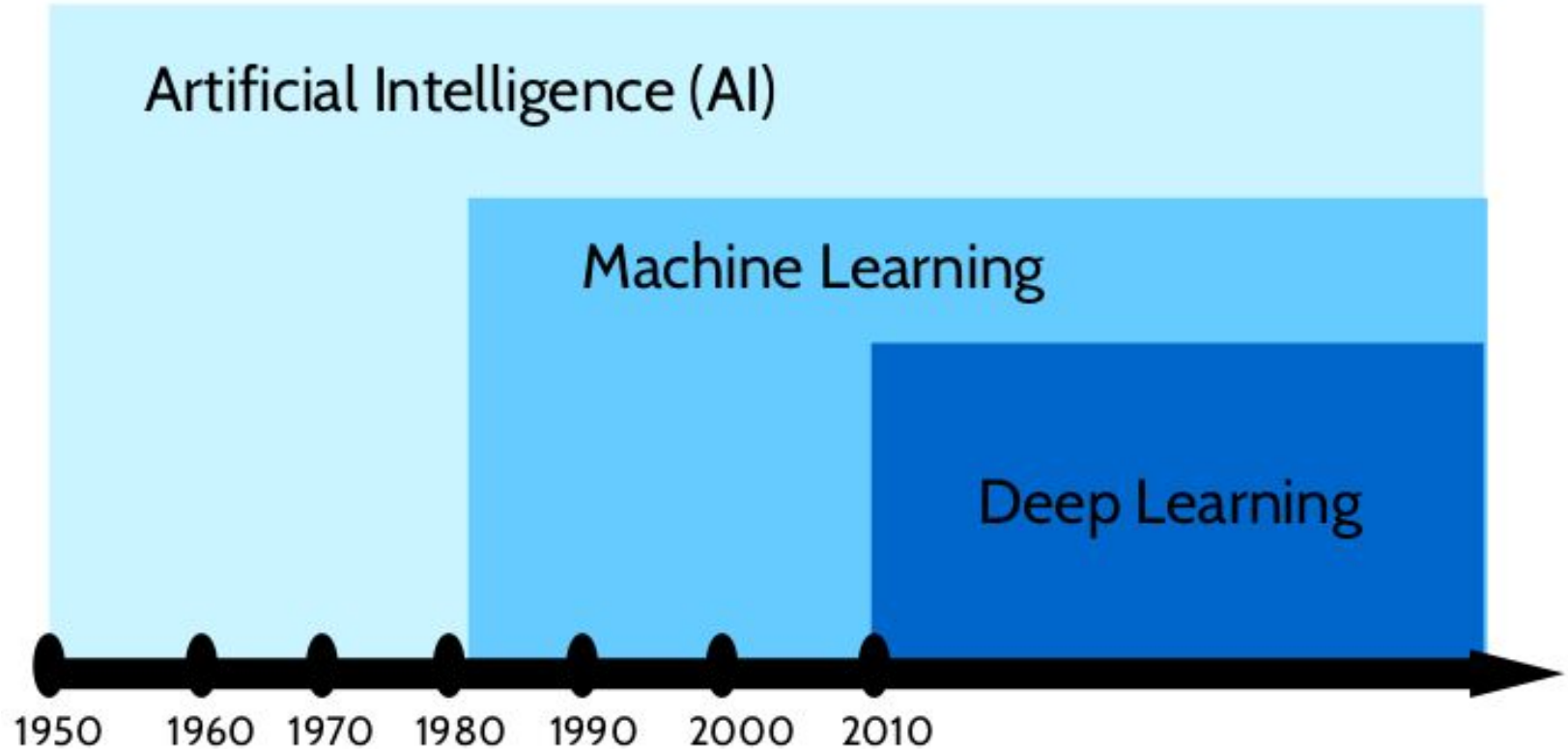
- Big data
- Analytics
- ML

# Today

- **Deep learning**
  - Neural networks: brief history, main architectures and how to train them



*Deep Learning is a subfield of ML concerned with algorithms inspired by the structure and function of the brain called artificial neural networks* [Jason Brownlee]

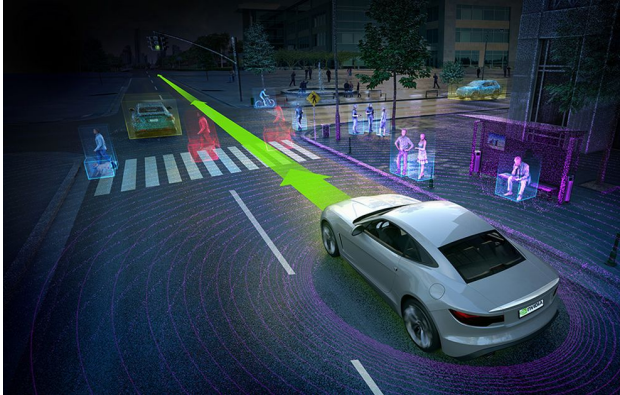
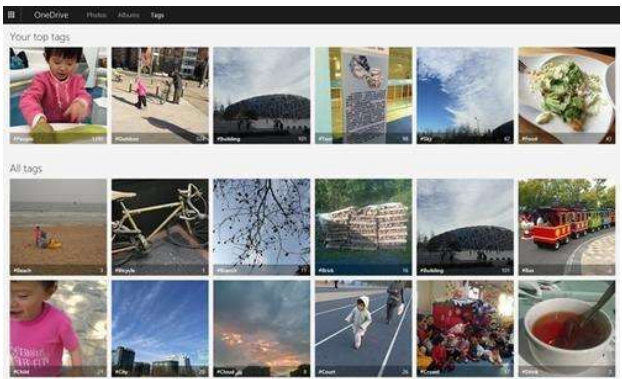


# Machine translation

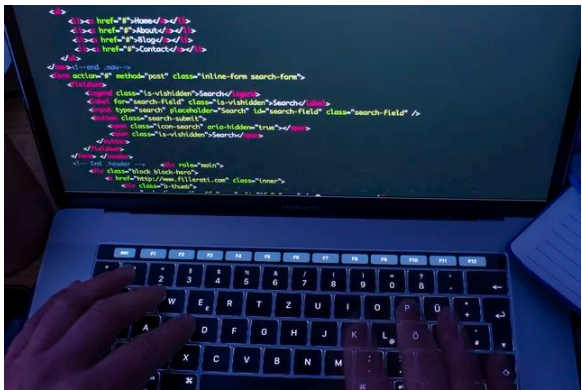
Real-time translation into  
Mandarin Chinese (2012)



# Visual recognition, CNN better than humans (2012)



Self driving cars, Tesla  
autopilot (2014)



Chatbots, CHAT GPT (2022)

# Creativity

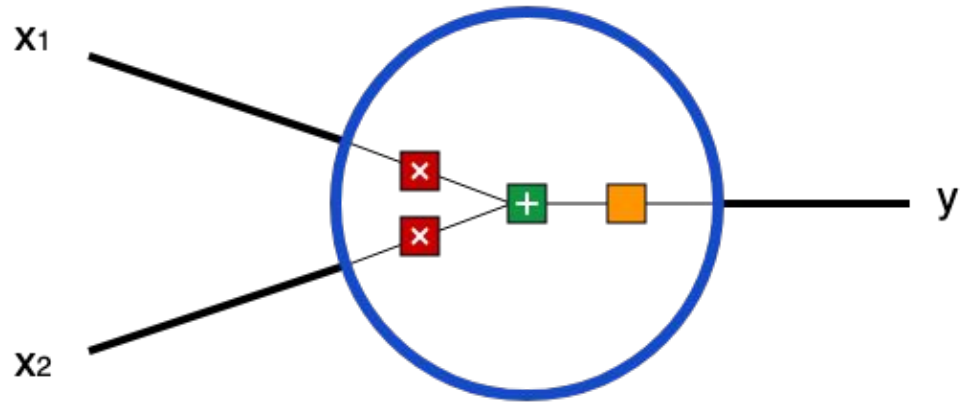
NST (2015), GAN (2014)



# Strategy games DeepMind beats Go world champion (2017)



# Short recap: Neuron



- each **input**  $x$  is multiplied by a **weight**  $w$

$$x_1 \rightarrow x_1 * w_1$$



$$x_2 \rightarrow x_2 * w_2$$

- all the weighted inputs are added together with a **bias**  $b$

$$(x_1 * w_1) + (x_2 * w_2) + b$$



- the sum is passed through an **activation function**  $f$

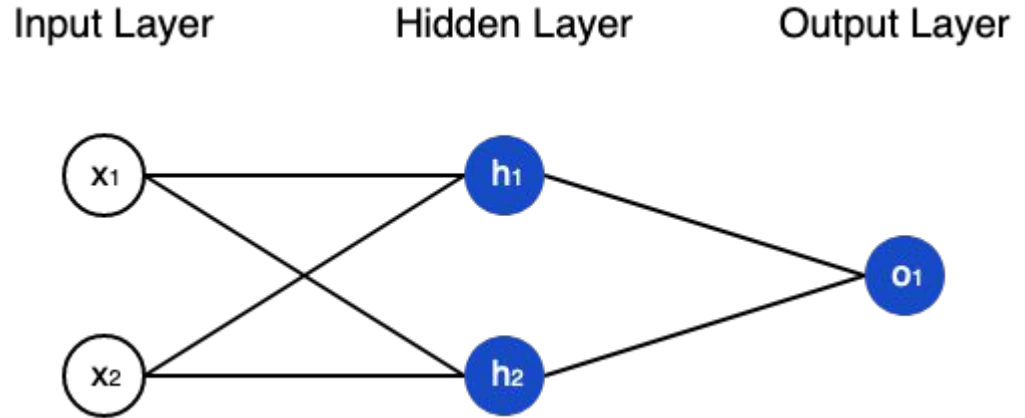
$$y = f(x_1 * w_1 + x_2 * w_2 + b)$$





# Neural network

- Combining more neurons
- A **hidden layer** is any layer between the input (first) layer and output (last) layer
  - There can be multiple hidden layers
- **Feedforward**: process of passing inputs forward to get an output

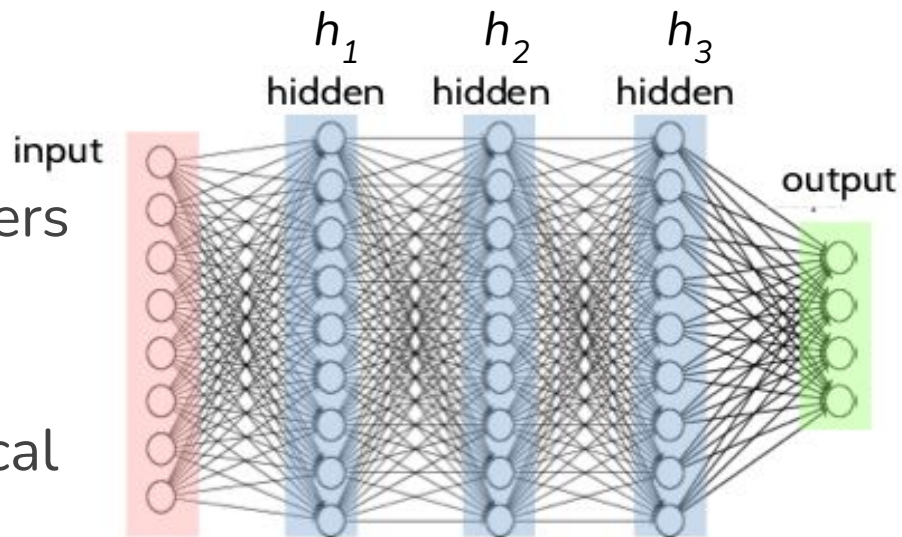


This network has:

- one **input** layer with 2 inputs
- one **hidden** layer with 2 neurons
- one **output** layer with 1 neuron

# Deep Learning

- Neural network with several layers
  - Deep vs shallow
- A family of parametric models which learn non-linear hierarchical representations:

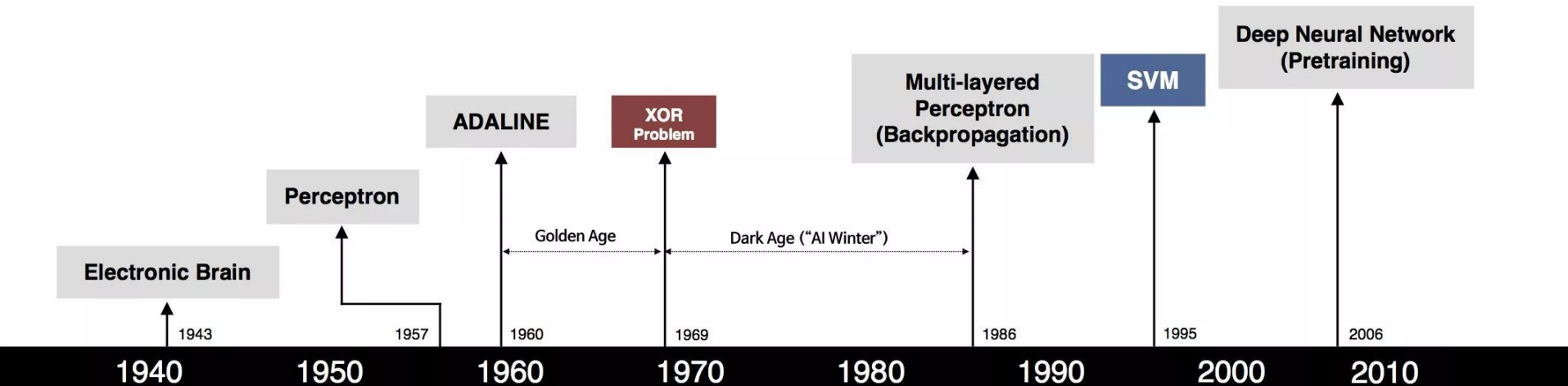


$$a_L(\mathbf{x}; \Theta) = h_L(h_{L-1}(\dots(h_1(\mathbf{x}, \theta_1), \theta_{L-1}), \theta_L)$$

Diagram illustrating the mathematical representation of a deep neural network's output function. The equation is shown with arrows pointing to specific components:

- $\mathbf{x}$ : input
- $\Theta$ : parameters of the network
- $h$ : non-linear activation function
- $\theta_L$ : parameters of layer  $L$

# Brief history of neural networks



S. McCulloch – W. Pitts



F. Rosenblatt



B. Widrow – M. Hoff



M. Minsky – S. Papert



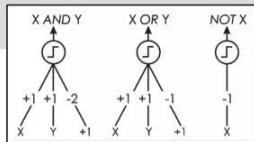
D. Rumelhart – G. Hinton – R. Williams



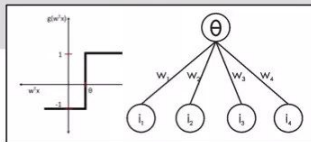
V. Vapnik – C. Cortes



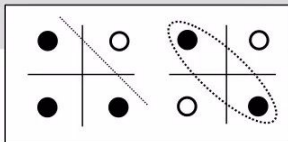
G. Hinton – S. Ruslan



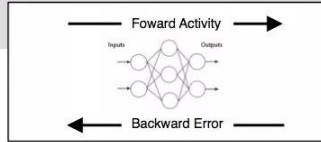
- Adjustable Weights
- Weights are not Learned



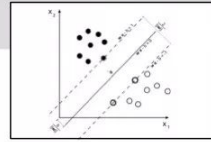
- Learnable Weights and Threshold



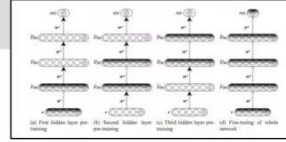
- XOR Problem



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



- Limitations of learning prior knowledge
- Kernel function: Human Intervention

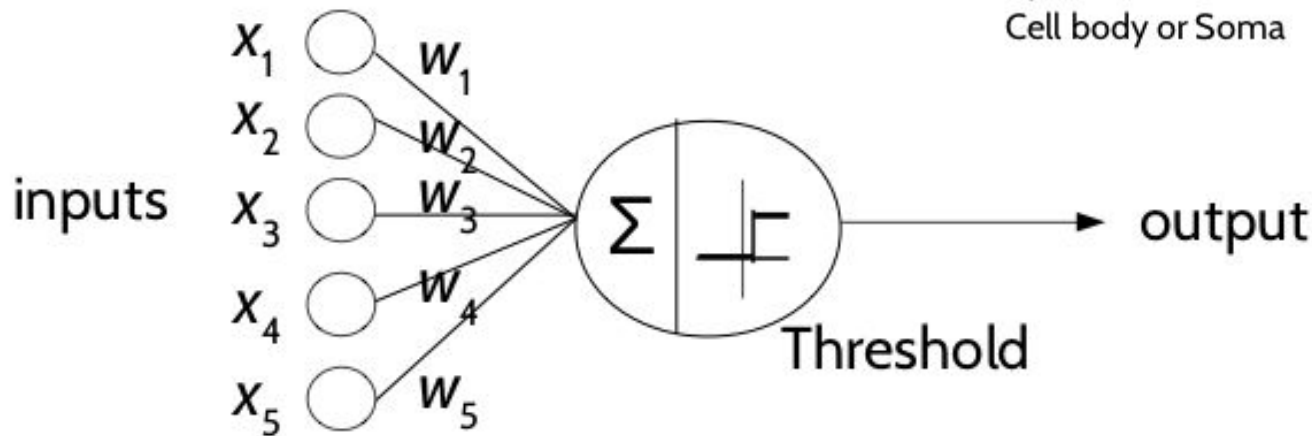
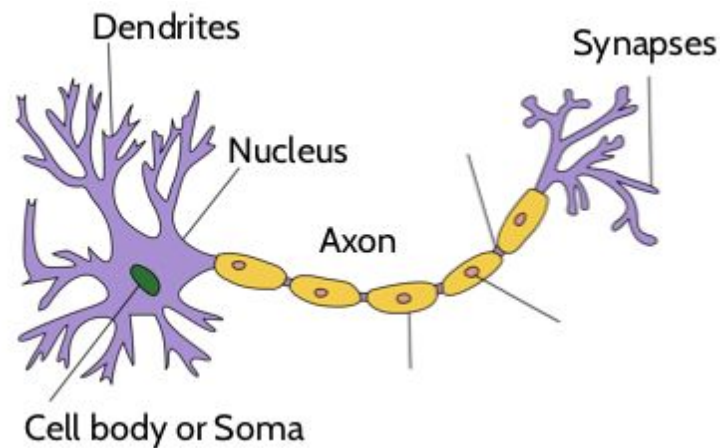


- Hierarchical feature Learning



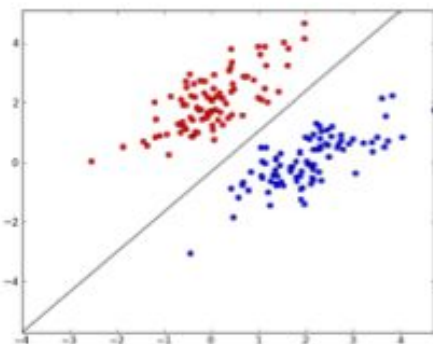
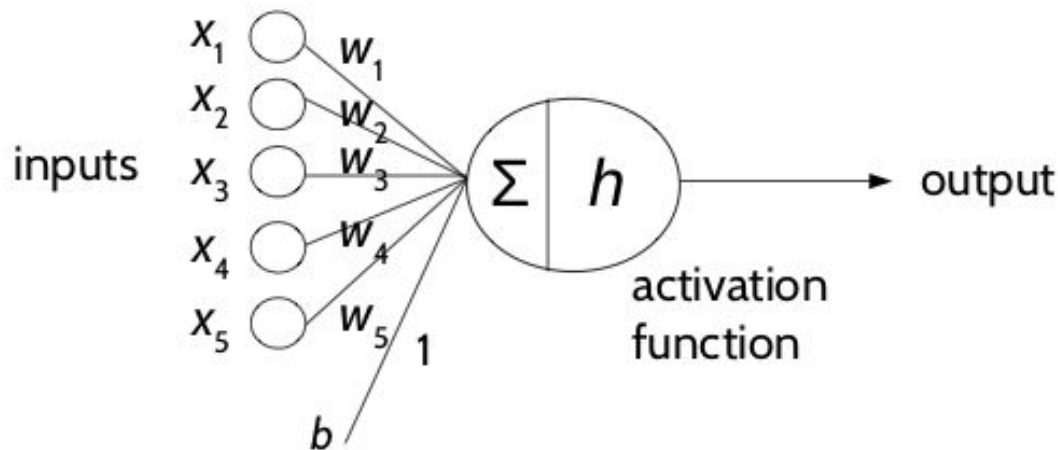
# 1943 – McCulloch & Pitts Model

- Early model of artificial neuron
- Generates a binary output
- The weights values are fixed



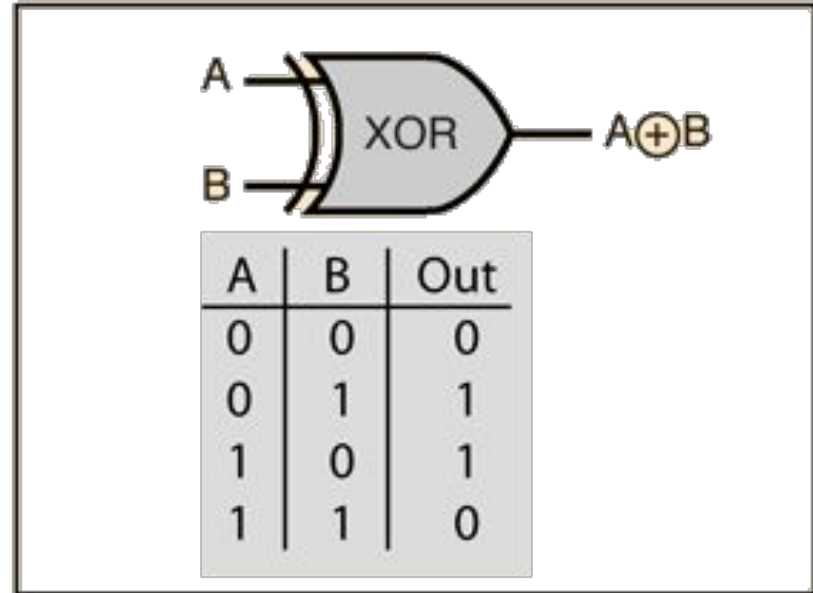
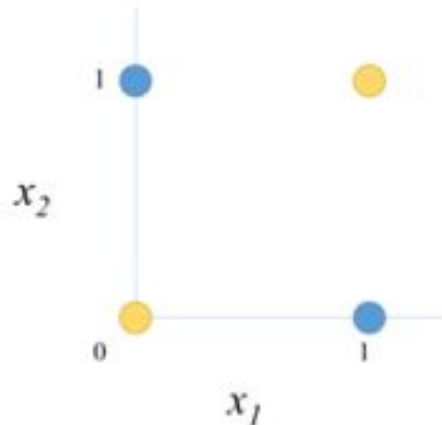
# 1958 – Perceptron by Rosemblatt

- Perceptron as a machine for linear classification
- Main idea: Learn the weights and consider bias.
  - One weight per input
  - Multiply weights with respective inputs and add bias
  - If result larger than **threshold** return 1, otherwise 0



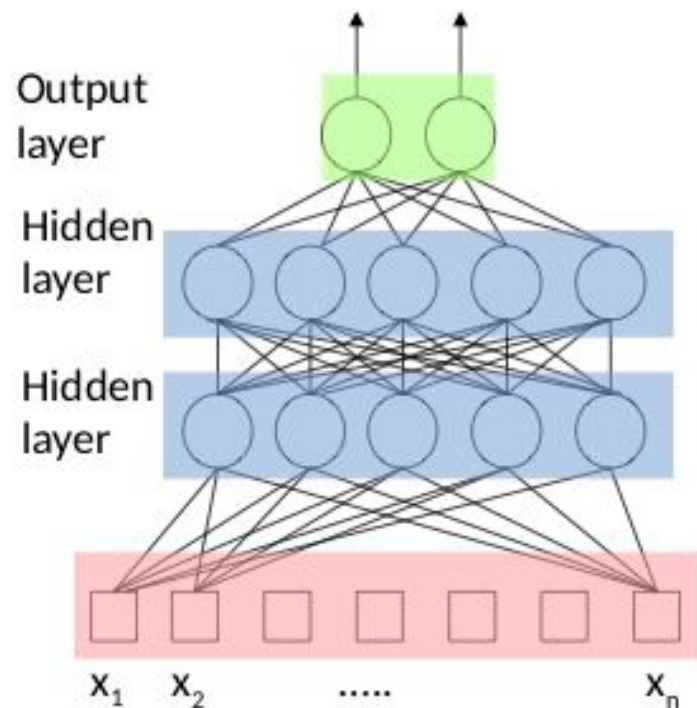
# First NN winter

- 1970- Minsky. The XOR cannot be solved by perceptrons.
- Neural models cannot be applied to complex tasks.



# Multi-layer Feed Forward Neural Network

- 1980s. Multi-layer Perceptrons (MLP) can solve XOR.
- ML Feed Forward Neural Networks:
  - Densely connect artificial neurons to realize compositions of non-linear functions
  - The information is propagated from the inputs to the outputs
  - The input data are usually  $n$ -dimensional feature vectors
  - Tasks: Classification, Regression



# How to train it?

- Rosenblatt algorithm\* not applicable, as it expects to know the desired target
  - For hidden layers we cannot know the desired target
- Learning MLP for complicated functions can be solved with **Back propagation\*\* (1980)**
  - efficient algorithm for complex NN which processes large training sets

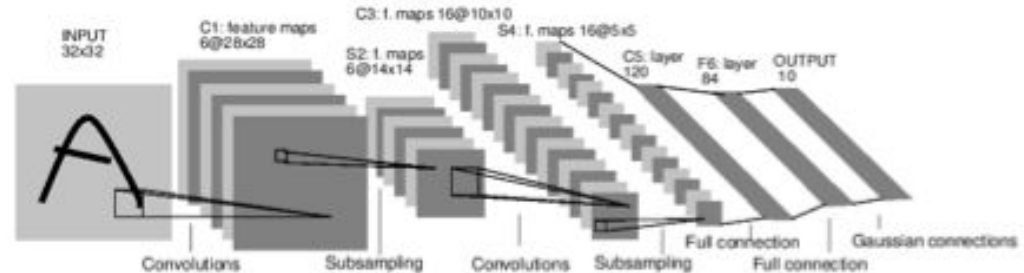
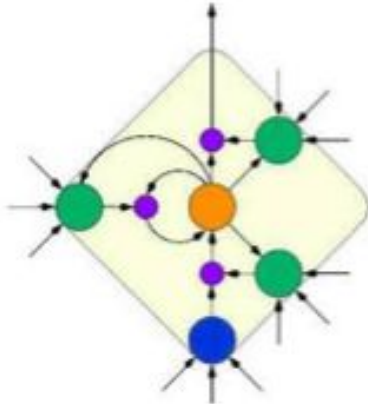
\* Remember? Rosenblatt developed a method to train a single neuron

\*\* later today



# 1990s - CNN and LSTM

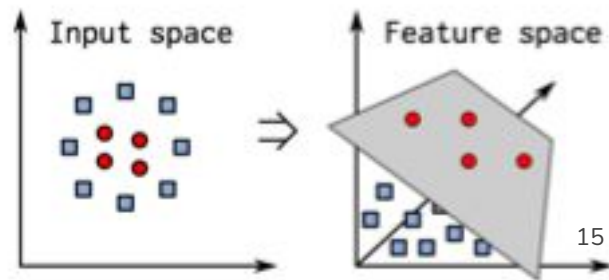
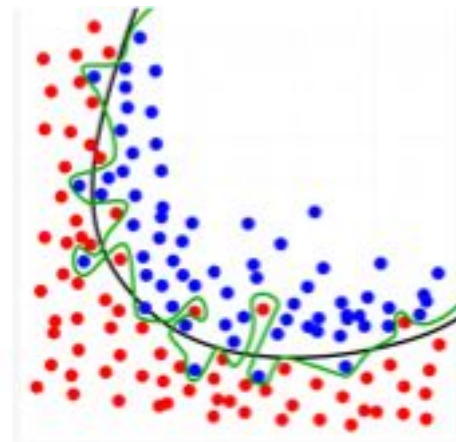
- Important advances in the field:
  - Backpropagation
  - Recurrent Long-Short Term Memory Networks (Schmidhuber, 1997)
  - Convolutional Neural Networks - LeNet: OCR solved before 2000s (LeCun, 1998).



OCR: Optical character recognition

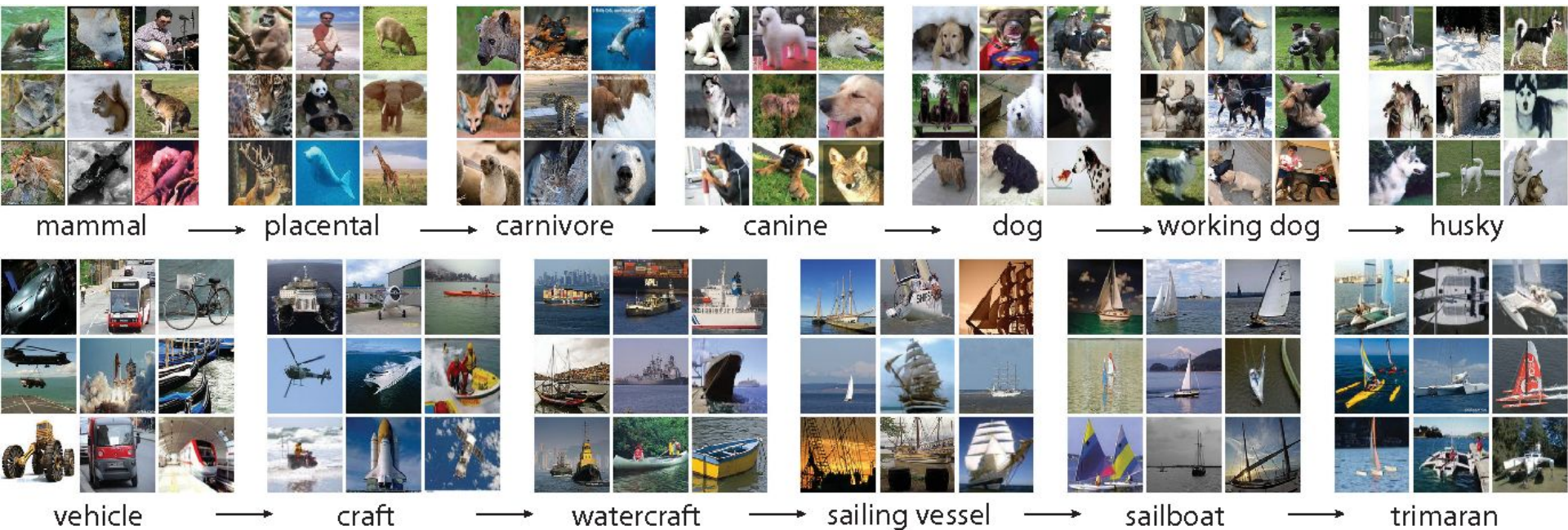
# Second NN Winter

- NN cannot exploit many layers
  - Overfitting
  - Vanishing gradient (with NN training you need to multiply several small numbers → they become smaller and smaller)
- Lack of processing power (no GPUs)
- Lack of data (no large annotated datasets)
- Kernel Machines (e.g. SVMs) suddenly become very popular◦



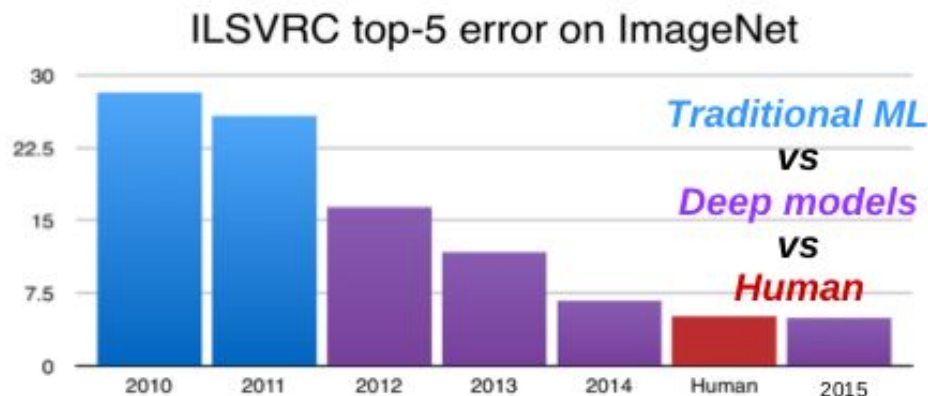
# ImageNet

A Large-Scale Hierarchical Image Database (2009)



# 2012 - AlexNet

- Hinton's group implemented a CNN similar to LeNet [LeCun1998] but...
  - Trained on ImageNet (1.4M images, 1K categories)
  - With 2 GPUs
  - Other technical improvements (ReLU, dropout, data augmentation)





# Convolutional Neural Networks (CNN)

- **Convolutional** layer: two functions produce a third that describes how the shape of one is changed by the other
- **pooling** layer: reduce dimensionality

Source layer

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

Convolutional kernel

-1	0	1
2	1	2
1	-2	0

Destination layer

		5					

$$\begin{aligned} &(-1 \times 5) + (0 \times 2) + (1 \times 6) + \\ &(2 \times 4) + (1 \times 3) + (2 \times 4) + \\ &(1 \times 3) + (-2 \times 9) + (0 \times 2) = 5 \end{aligned}$$

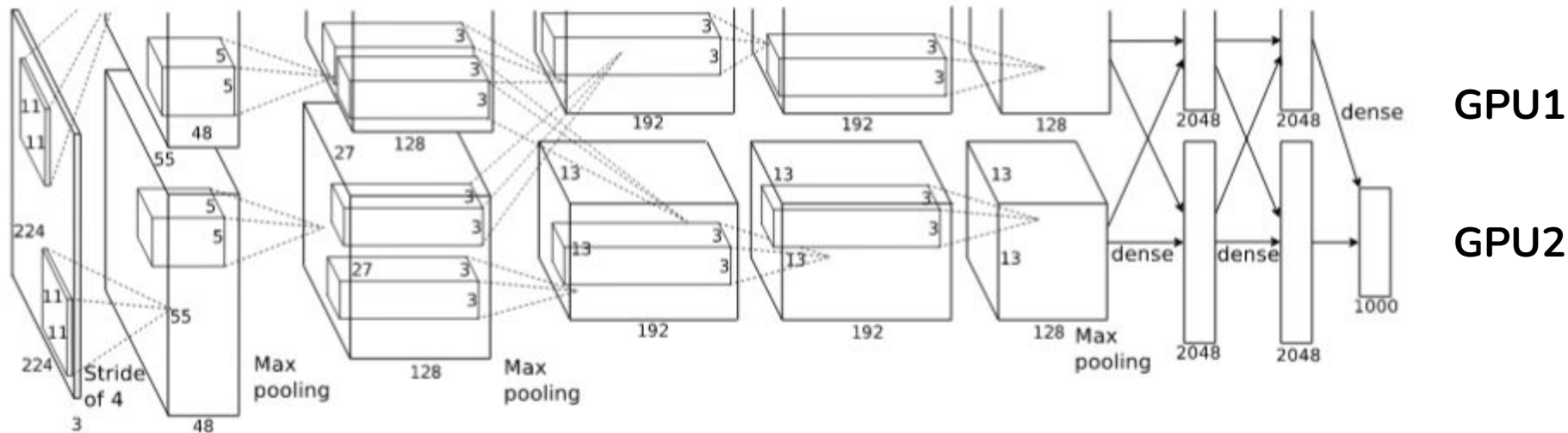
4	1	5	0
7	8	9	8
3	5	6	5
2	4	1	0

Max pooling

8	9
5	6



# AlexNet

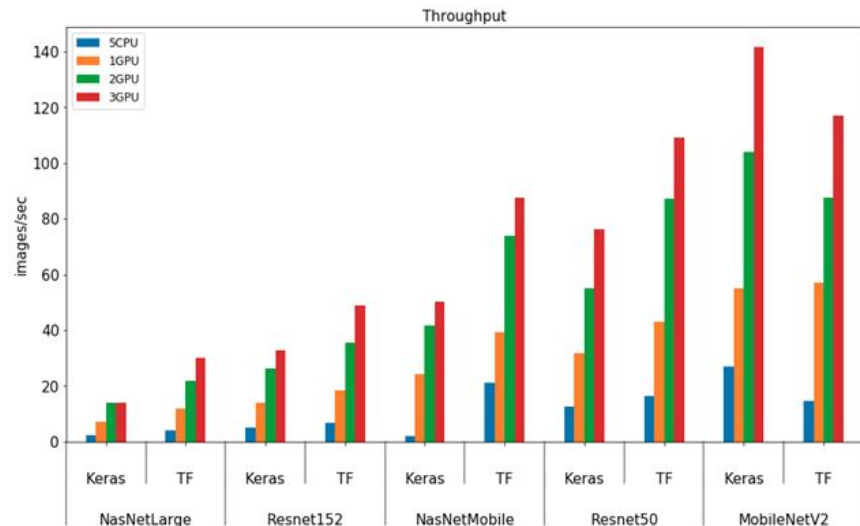


- 60M parameters
- Limited information exchange between GPUs

# Why Deep Learning now?

- Three main factors:
  - Better hardware
  - Big data
- Technical advances:
  - Layer-wise pretraining
  - Optimization (e.g. Adam, batch normalization)
  - Regularization (e.g. dropout)

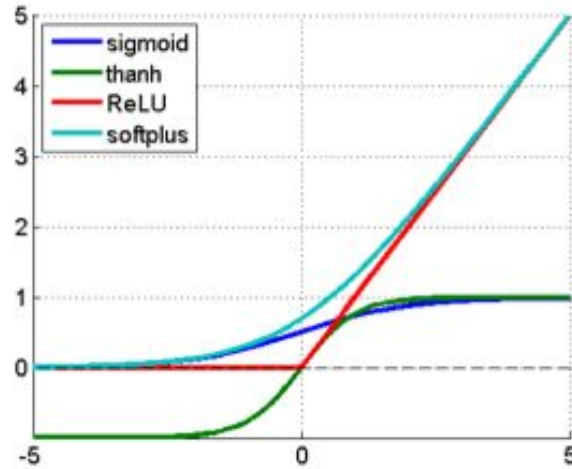
....



# Rectified Linear Units - Activation function (2010)

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

- More efficient gradient propagation: (derivative is 0 or constant)
- More efficient computation: (only comparison, addition and multiplication).
- Sparse activation: e.g. in a randomly initialized networks, only about 50% of hidden units are activated (having a non-zero output)

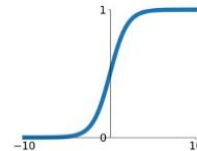


# Activation functions

- Classification: sigmoid functions
  - sigmoids and tanh functions are sometimes avoided due to the vanishing gradient problem
- **ReLU** function is a general activation function
- dead neurons in our networks -> the leaky ReLU
- ReLU function should only be used in the hidden layers
- As a rule of thumb, start with ReLU

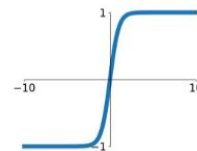
## Sigmoid

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



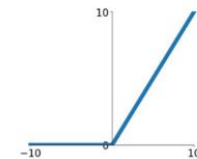
## tanh

$$\tanh(x)$$



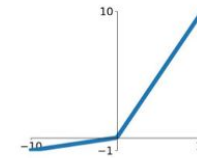
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

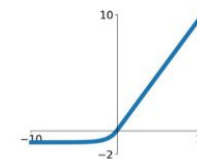


## Maxout

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

## ELU

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



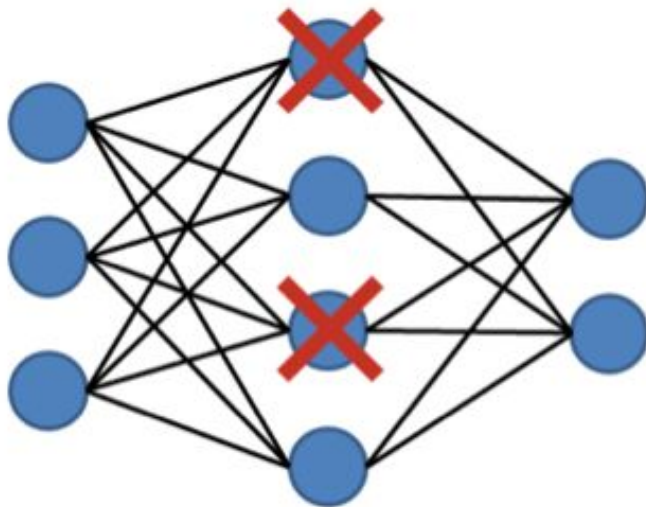
# Regularization

- One of the major aspects of training the model is overfitting -> the ML model captures the noise in your training dataset
- The **regularization** term is an addition to the loss function which helps generalize the model
  - **L1** or Lasso regularization adds a penalty which is the sum of the absolute values of the weights
    - L1+MSE
$$\text{Min}(\sum_{i=1}^n (y_i - w_i x_i)^2 + p \sum_{i=1}^n |w_i|)$$
  - **L2** or Ridge regularization adds a penalty which is the sum of the squared values of weights
    - L2+MSE
$$\text{Min}(\sum_{i=1}^n (y_i - w_i x_i)^2 + p \sum_{i=1}^n (w_i)^2)$$
- **Early Stopping** is a time regularization technique which stops training based on given criteria



# Regularization - Dropout

- For each instance drop a node (hidden or input) and its connections with probability  $p$  and train
- Final net just has all averaged weights (actually scaled by  $1-p$ )
- As if ensembling  $2^n$  different network substructures



# Data augmentation

- Techniques to significantly increase the diversity of data available for training models, without actually collecting new data
- Data augmentation techniques such as cropping, padding, and horizontal flipping are commonly used to train large neural networks



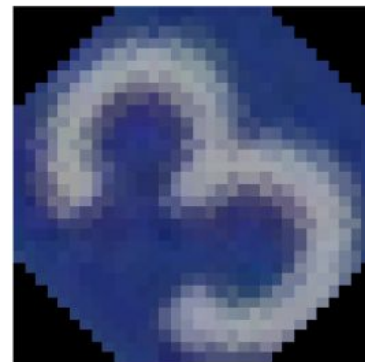
Original



Horizontal Flip

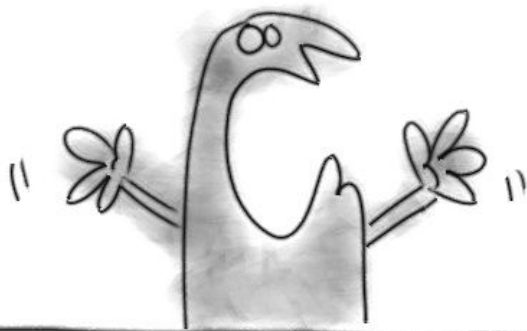


Pad & Crop



Rotate

# Now What?!!



# Training a neural network

- **Problem:** Predict gender from weight and height
- **Input Dataset:**

Name	Features		Labels
	Weight (lb)	Height (in)	Gender
Alice	133	65	F
Bob	160	72	M
Charlie	152	70	M
Diana	120	60	F

# 1. Feature engineering

- Symmetrize numeric values
- Category -> numbers

Name	Weight (lb)	Height (in)	Gender
Alice	133	65	F
Bob	160	72	M
Charlie	152	70	M
Diana	120	60	F

Name	Weight (minus 135)	Height (minus 66)	Gender
Alice	-2	-1	1
Bob	25	6	0
Charlie	17	4	0
Diana	-15	-6	1





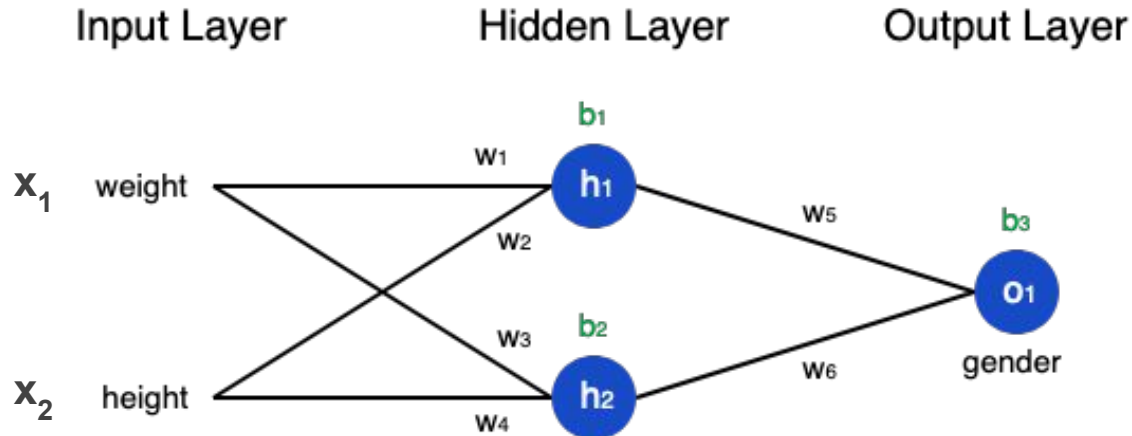
# Ingredients

- $n$  : 4, number of samples (Alice, Bob, Charlie, Diana)
- **Inputs:**  $\mathbf{X}$ , dimension = 2
  - **weights** ( $\mathbf{X}_1$ ) and **heights** ( $\mathbf{X}_2$ )
- **Features:**  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ , transformed inputs
- $\mathbf{y}$  : variable being predicted (Gender)
- $\mathbf{y}_{\text{true}}$  : true value of  $\mathbf{y}$

Name	Weight (minus 135)	Height (minus 66)	Gender
Alice	-2	-1	1
Bob	25	6	0
Charlie	17	4	0
Diana	-15	-6	1

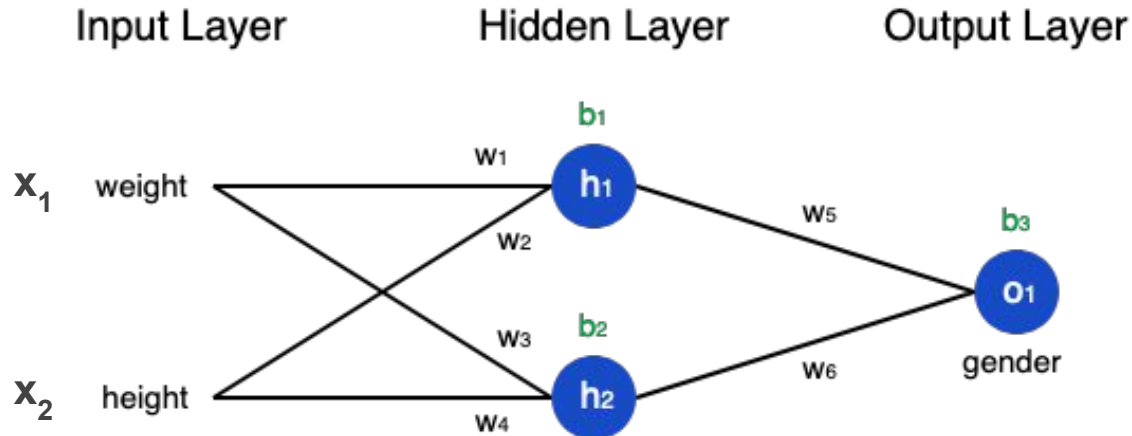
## 2. Model

- $y_{\text{pred}}$  : predicted value of  $y = \text{gender}$
- Neural network: with 1 hidden layer with 2 neurons
- Outputs of the hidden layer:  $\mathbf{h}$
- Unknown parameters: weights  $\mathbf{w}$  and biases  $\mathbf{b}$



# Some math

- For each **neuron**:  $y_j = b_j + \mathbf{f} \sum_i x_i w_{ij}$ ,  $\mathbf{f}$  is the activation function
- For the **net**:
  - $h_1 = \mathbf{f}(w_1 x_1 + w_2 x_2) + b_1$
  - $h_2 = \mathbf{f}(w_3 x_1 + w_4 x_2) + b_2$
  - $o_1 = \mathbf{f}(w_5 h_1 + w_6 h_2) + b_3$



### 3. Model training

- Training the network == Find weights **w** and biases **b** that minimize the loss
- Loss function **L**: **MSE**
- Find weights **w** and biases **b** to minimise

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_{true} - y_{pred})^2$$

$$L(w_1, w_2, w_3, w_4, w_5, w_6, b_1, b_2, b_3)$$

- Typically, this is already implemented in ML software packages

# Back propagation

- Minimization taking **partial derivatives**
- For very simple case: with only Alice in the dataset,  $n=1$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y_{pred}} * \frac{\partial y_{pred}}{\partial h_1} * \frac{\partial h_1}{\partial w_1}$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_{true} - y_{pred})^2$$

$$L = (1 - y_{pred})^2$$

$$\frac{\partial L}{\partial y_{pred}} = \frac{\partial (1 - y_{pred})^2}{\partial y_{pred}}$$

$$y_{pred} = o_1 = f(w_5 h_1 + w_6 h_2 + b_3)$$

$$\frac{\partial y_{pred}}{\partial h_1} = w_5 * f'(w_5 h_1 + w_6 h_2 + b_3)$$

$$h_1 = f(w_1 x_1 + w_2 x_2 + b_1)$$

$$\frac{\partial h_1}{\partial w_1} = x_1 * f'(w_1 x_1 + w_2 x_2 + b_1)$$

# Stochastic Gradient Descent (SGD)

- Start with randomly initialised weights

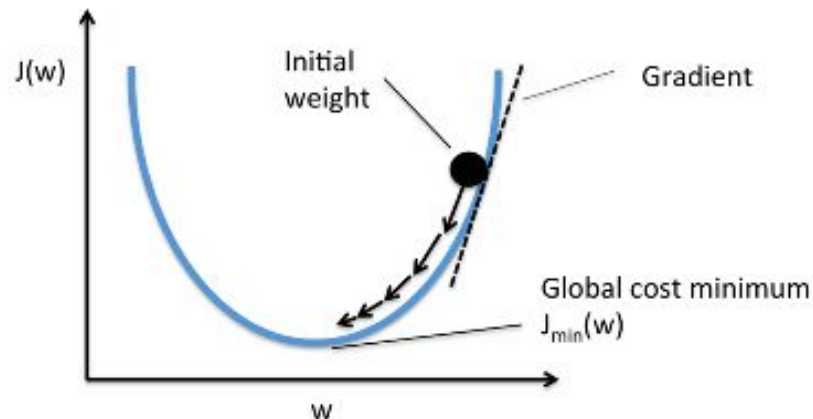
- **update equation:**

$$w_1 \leftarrow w_1 - \eta \frac{\partial L}{\partial w_1}$$

- $\eta$  is a constant called the **learning rate** that controls how fast we train

- If  $\frac{\partial L}{\partial w_1}$  is positive,  $w_1$  will decrease, which makes  $L$  decrease.

- If  $\frac{\partial L}{\partial w_1}$  is negative,  $w_1$  will increase, which makes  $L$  decrease.



- **Stochastic vs Batch (BGD)** -> the parameters are updated using only one single training instance (usually randomly selected) in each iteration vs the whole training set (== batch)



# Mini-batch gradient descent

- Use mini-batch **sampled** in the dataset for gradient estimate.
- Sometimes helps to escape from local minima
- Noisy gradients act as regularization
- Variance of gradients increases when batch size decreases
- Not clear how many sample per batch

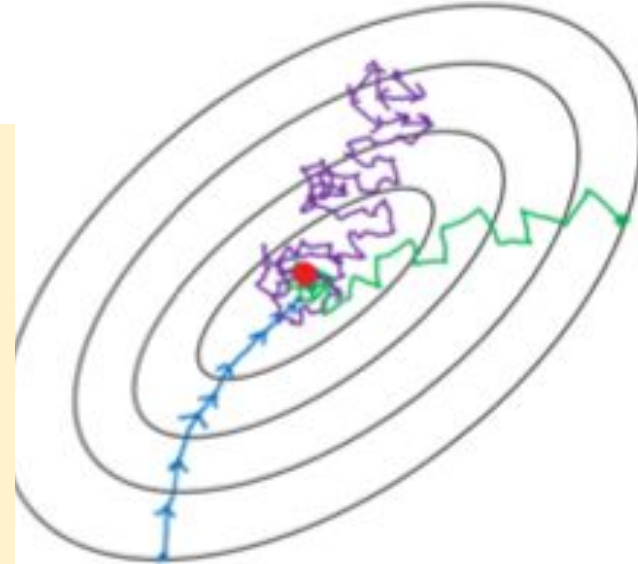
What happens in **one epoch** for SGD:

1. Take a random sample
2. Feed it to Neural Network
3. Calculate its gradient
4. Use the gradient we calculated in step 3 to update the weights
5. Repeat steps 1–4 for all the examples in training dataset

What happens in **one epoch** for Mini-BGD:

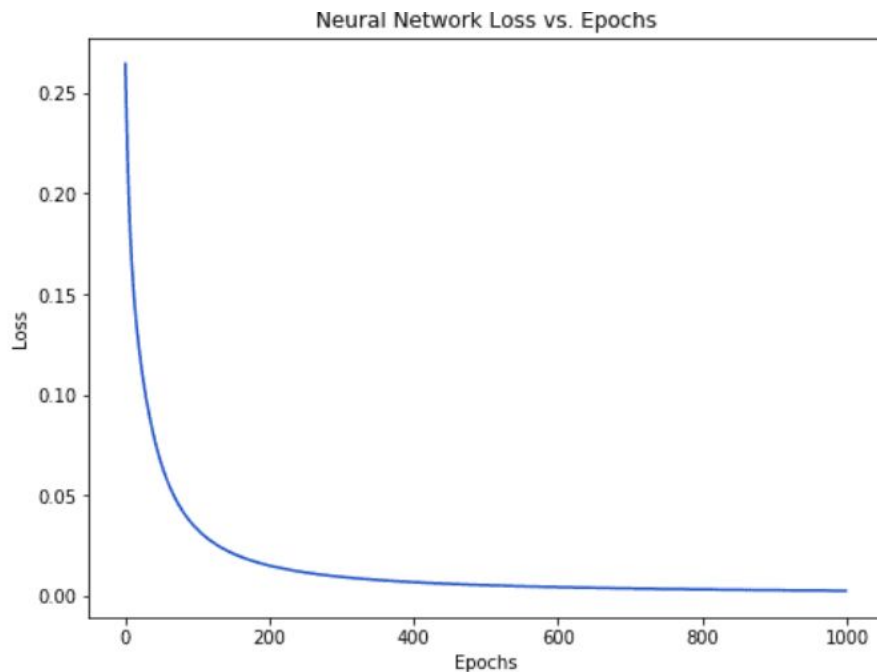
1. Pick a mini-batch
2. Feed it to Neural Network
3. Calculate the mean gradient of the mini-batch
4. Use the mean gradient we calculated in step 3 to update the weights
5. Repeat steps 1–4 for the mini-batches we created

— Batch gradient descent  
— Mini-batch gradient Descent  
— Stochastic gradient descent



# Take home messages

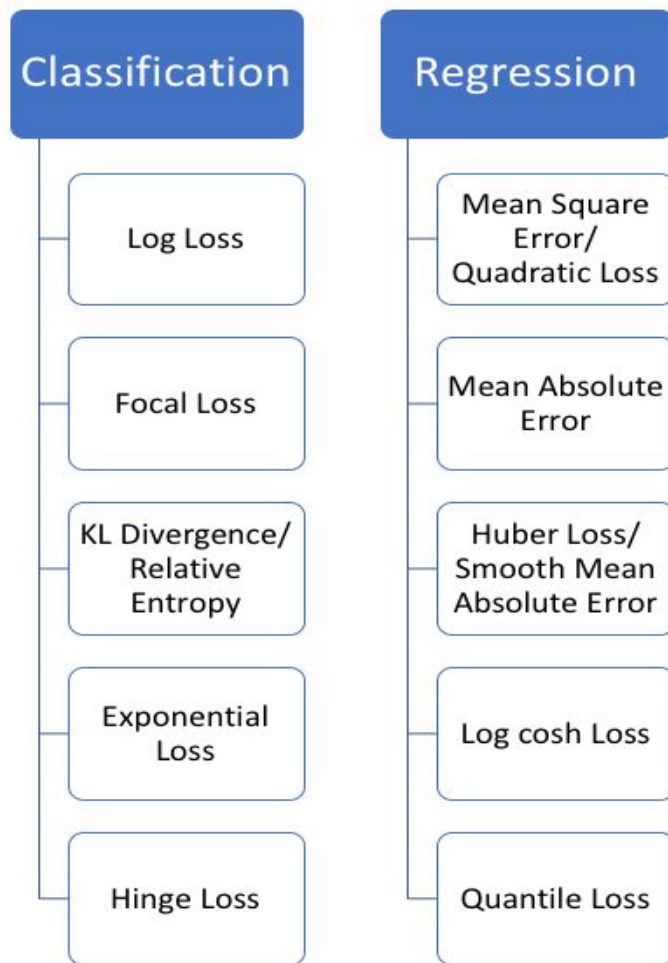
- **Gradient descent** is an iterative learning algorithm that uses a training dataset to update a model.
  - **BGD**: Batch Size = Size of Training Set
  - **SGD**: Batch Size = 1
  - **Mini-BGD**.  $1 < \text{Batch Size} < \text{Size of Training Set}$
- The **batch size** is a hyperparameter of gradient descent that controls the number of training samples to work through before the model's internal parameters are updated.



- The number of **epochs** is a hyperparameter of gradient descent that controls the number of complete passes through the training dataset

# Loss functions

- <https://heartbeat.fritz.ai/5-regression-loss-functions-all-machine-learners-should-know-4fb140e9d4b0>
- [https://www.wikiwand.com/en/Loss\\_functions\\_for\\_classification](https://www.wikiwand.com/en/Loss_functions_for_classification)



# Hands-on today

1. **Point your browser to:** <https://yoga.to.infn.it>
2. **Open a terminal:**
  - `cd MLCourse-2223`
  - `git pull`
  - `cp Notebooks/Day3/* ../`
3. **From JupyterHub Home tab:**
  - start and run *ML\_MCP.ipynb*
  - Apache ML Library [MLLib](#)
    - Multilayer Perceptron Classifier (MPC)
  - [Keras](#)
    - Sequential model