# Big data science Day 3

F. Legger - INFN Torino
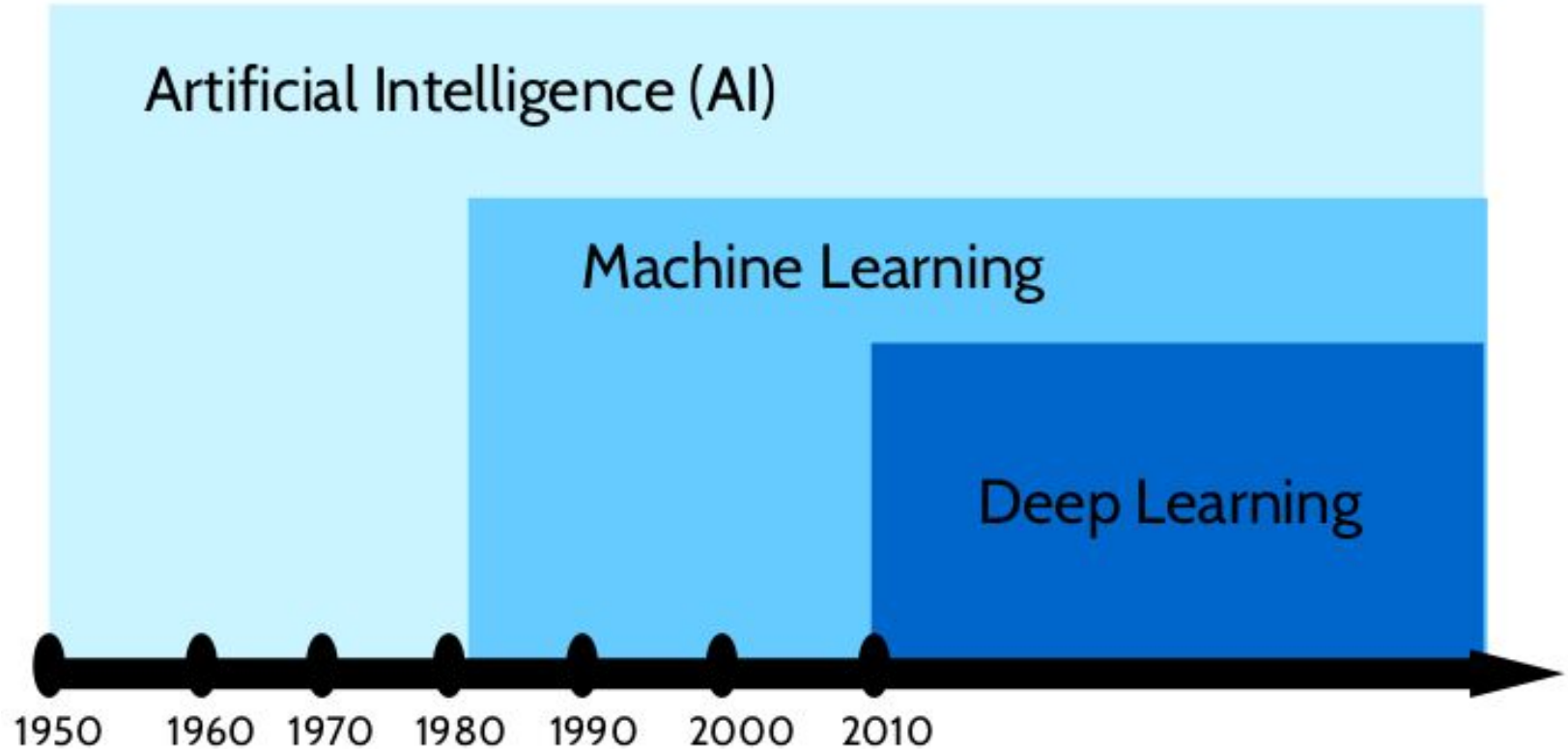https://github.com/leggerf/MLCourse-2022

# We learned

- Big data
- Analytics
- Machine learning

# Today

- **Deep learning**
- Parallelisation
- Heterogeneous architectures
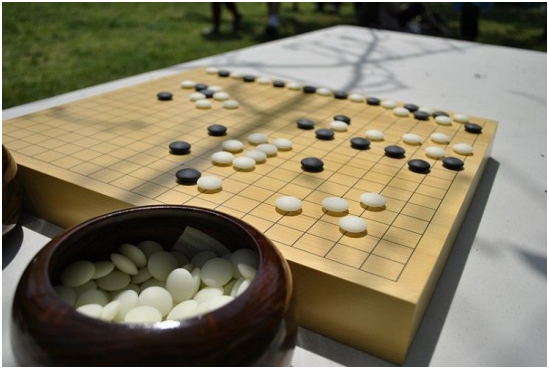- Future directions

*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]



Artificial Intelligence (AI)

Machine Learning

Deep Learning

1950    1960   1970   1980   1990   2000   2010

**Machine translation**
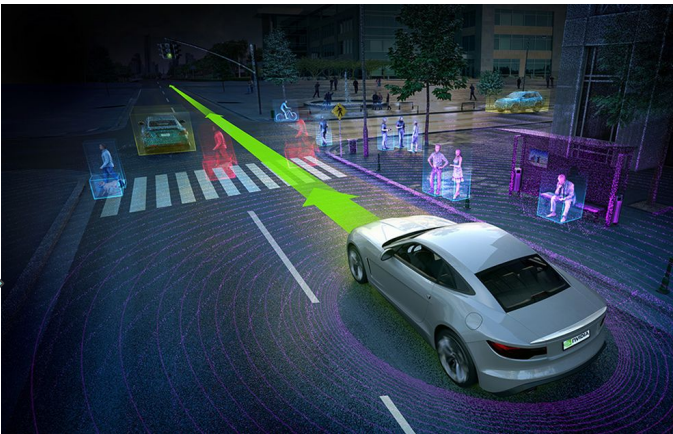Real-time translation into Mandarin Chinese (2012)

**Creativity**

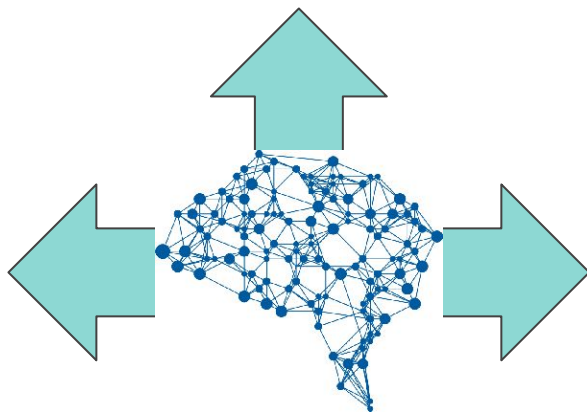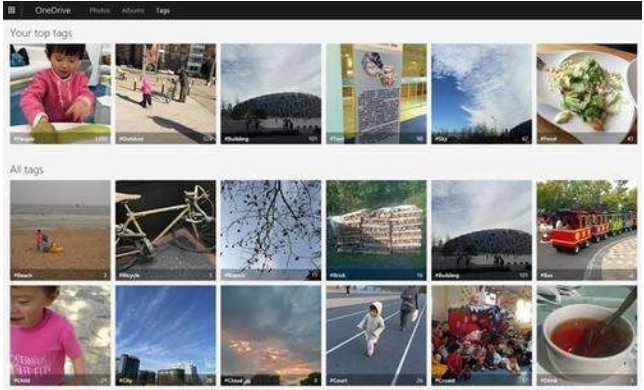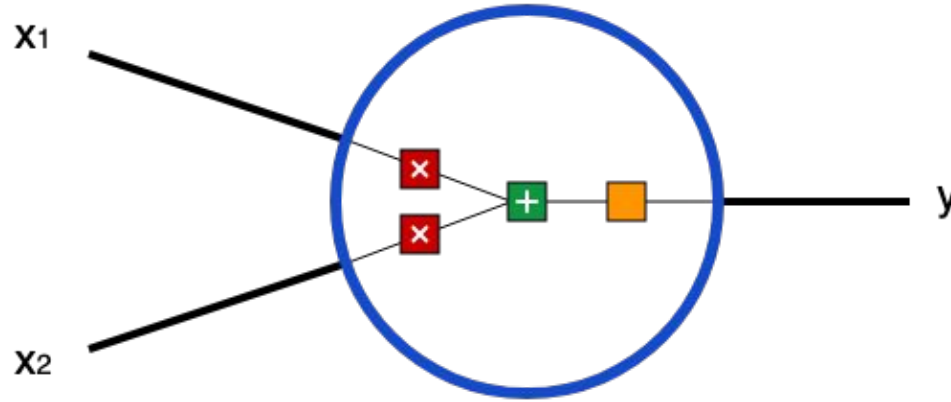**Strategy games**
DeepMind beats Go world champion (2017)

**Self driving cars**

**Visual recognition**

# Short recap: Neuron

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

$$x_1 \rightarrow x_1 * w_1$$

X

$$x_2 \rightarrow x_2 * w_2$$

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

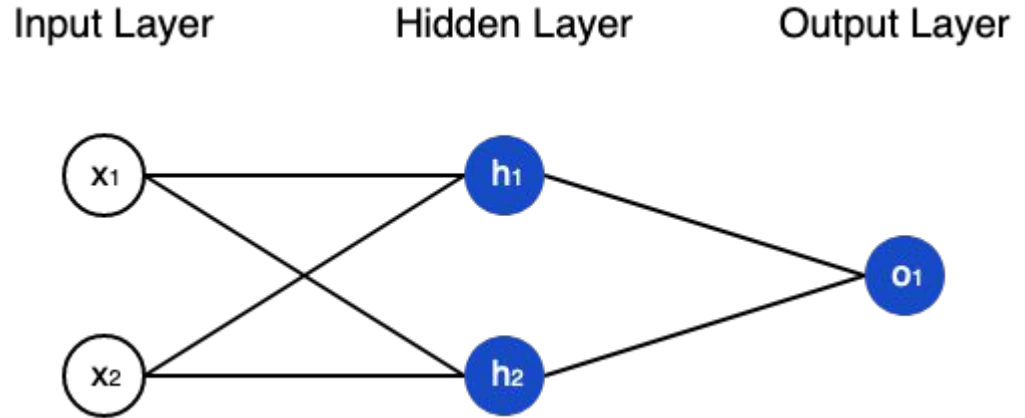$$(x_1 * w_1) + (x_2 * w_2) + b$$

+

- 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
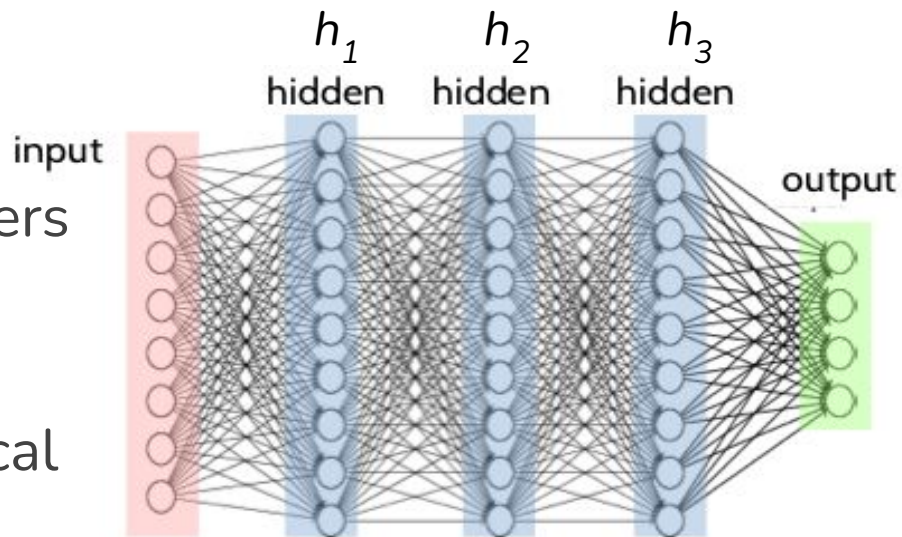


Input Layer     Hidden Layer     Output Layer

This network has:
  - one **input** layer with 2 inputs
  - one **hidden** layer with 2 neurons
  - one **output** layer with 1 neuron

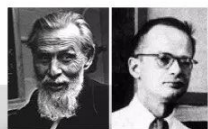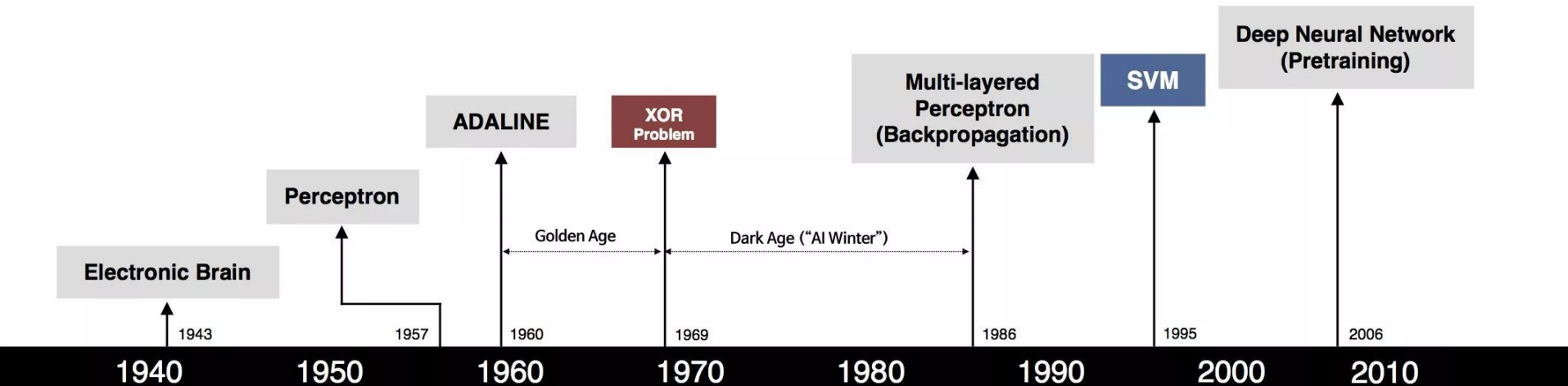# Deep Learning

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



$$a_L(\mathbf{x}; \boldsymbol{\Theta}) = h_L(h_{L-1}(...(h_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_{L-1}), \boldsymbol{\theta}_L)$$

input

parameters of the network

non-linear activation function

parameters of layer $L$

# Brief history of neural networks



Deep Neural Network (Pretraining)

SVM

Multi-layered Perceptron (Backpropagation)

ADALINE

XOR Problem

Perceptron

Golden Age

Dark Age ("AI Winter")

Electronic Brain

1943        1957        1960        1969                      1986              1995              2006

1940     1950     1960     1970     1980     1990     2000     2010

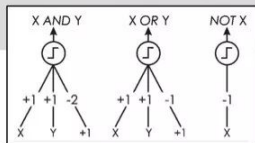S. McCulloch – W. Pitts     F. Rosenblatt     B. Widrow – M. Hoff     M. Minsky – S. Papert     D. Rumelhart – G. Hinton – R. Wiliams     V. Vapnik – C. Cortes     G. Hinton – S. Ruslan
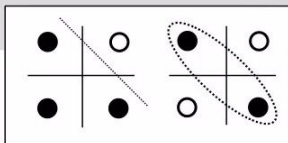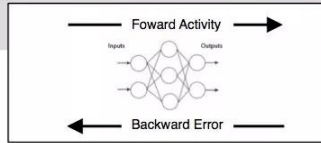
X AND Y     X OR Y     NOT X

+1 +1 -2     +1 +1 -1     -1

X Y     X Y     X

• Adjustable Weights
• Weights are not Learned

• Learnable Weights and Threshold

• XOR Problem

Foward Activity

Inputs     Outputs

Backward Error
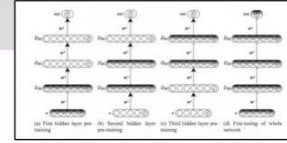
• Solution to nonlinearly separable problems
• Big computation, local optima and overfitting

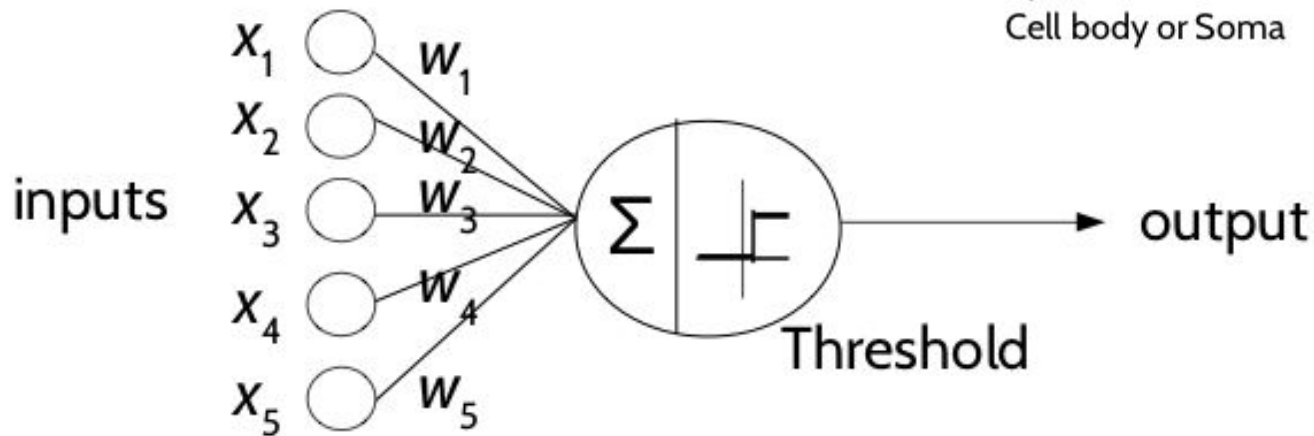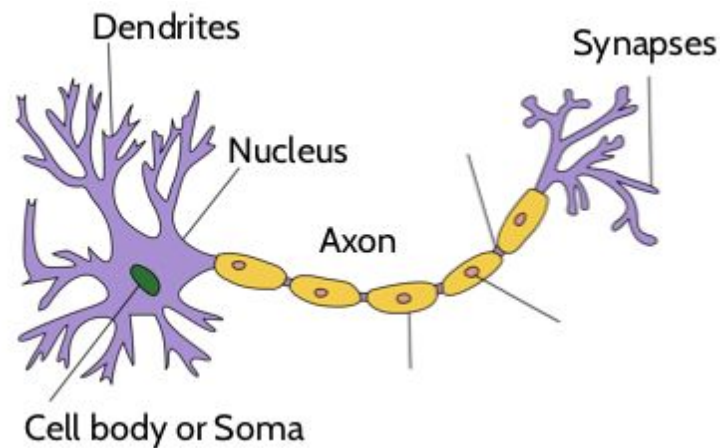• Limitations of learning prior knowledge
• Kernel function: Human Intervention

• Hierarchical feature Learning

8

# 1943 – McCulloch & Pitts Model

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

Dendrites

Synapses

Nucleus

Axon

Cell body or Soma

inputs

$x_1$  $w_1$
$x_2$  $w_2$
$x_3$  $w_3$
$x_4$  $w_4$
$x_5$  $w_5$

$\Sigma$

Threshold

output

# 1958 – Perceptron by Rosemblatt

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

inputs $x_1$ $w_1$ $x_2$ $w_2$ $x_3$ $w_3$ $\Sigma$ $h$ output $x_4$ $w_4$ $x_5$ $w_5$ $1$

activation function

$b$

# First NN winter

- 1970- Minsky. The XOR cannot be solved by perceptrons.

- Neural models cannot be applied to complex tasks.



A — XOR — A⊕B
B —

| A | B | Out |
|---|---|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$x_2$

$x_1$

0    1

# 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

Output layer

Hidden layer

Hidden layer

$x_1$  $x_2$  .....  $x_n$

12

# 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
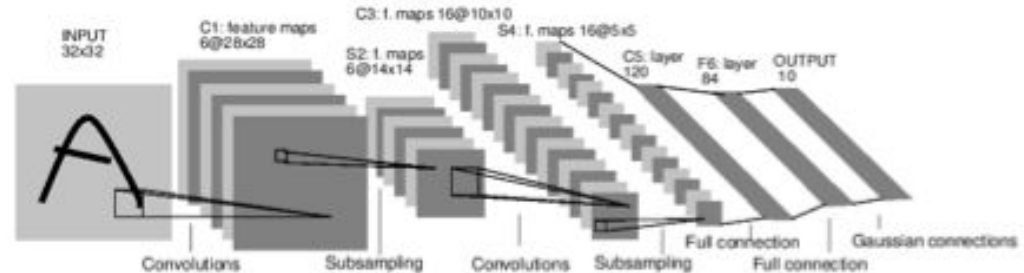- Learning MLP for complicated functions can be solved with **Back propagation (1980)**
  - efficient algorithm for complex NN which processes large training sets

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

# 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

# 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 | | | | | |
|---|---|---|---|---|---|
| | | | | | |

(-1×5) + (0×2) + (1×6) +
(2×4) + (1×3) + (2×4) +
(1×3) + (-2×9) + (0×2) = 5

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

Max pooling →

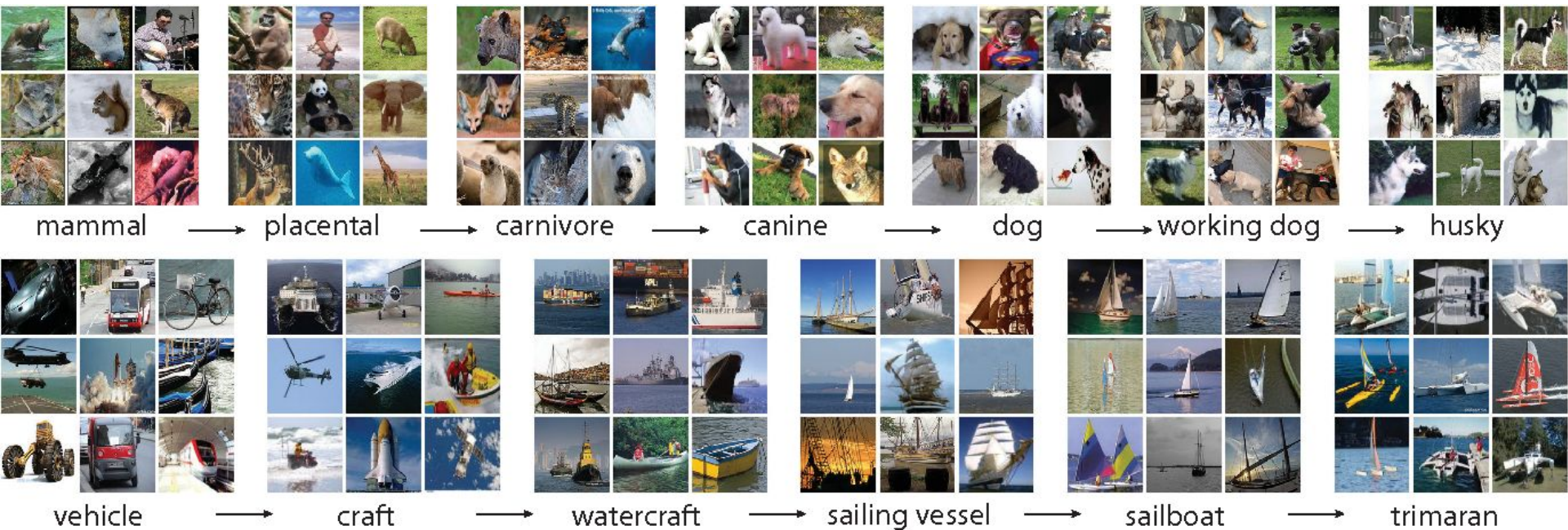| 8 | 9 |
|---|---|
| 5 | 6 |

# 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)
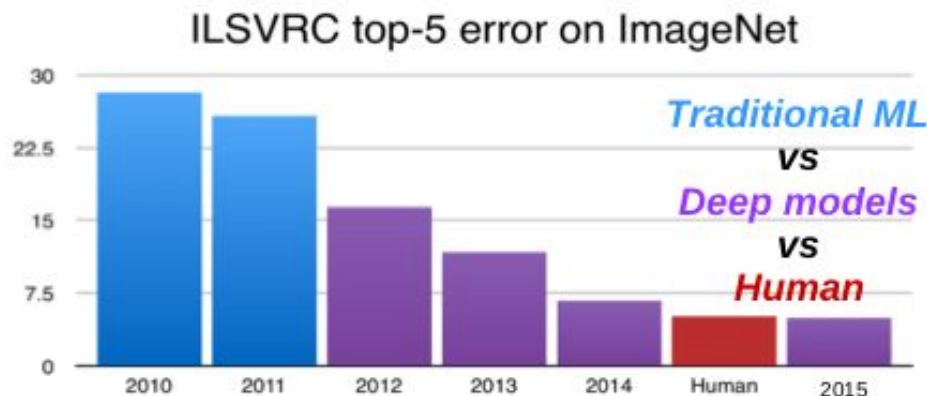- Kernel Machines (e.g. SVMs) suddenly become very popular

# ImageNet

A Large-Scale Hierarchical Image Database (2009)



mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

# 2012 – AlexNet

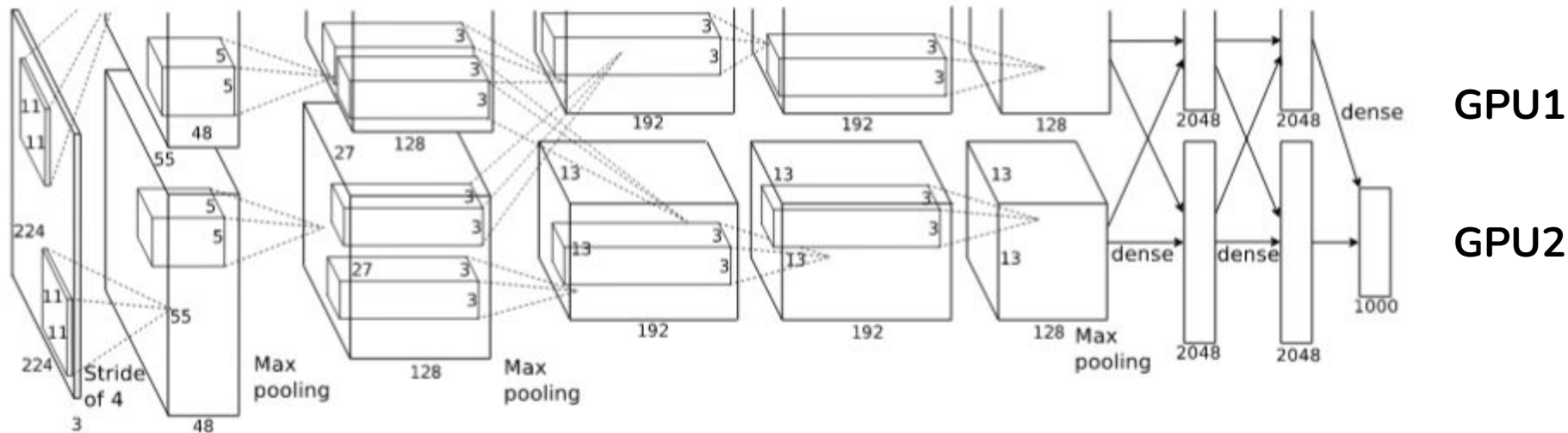- Hinton's group implemented a CNN similar to LeNet [LeCun1998] but...
  - Trained on ImageNet (1.4M images, 1K categories)
  - With 2 GPUs
  - Other technical improvements (ReLU, dropout, data augmentation)



ILSVRC top-5 error on ImageNet

**Traditional ML**
*vs*
**Deep models**
*vs*
**Human**

https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf 18
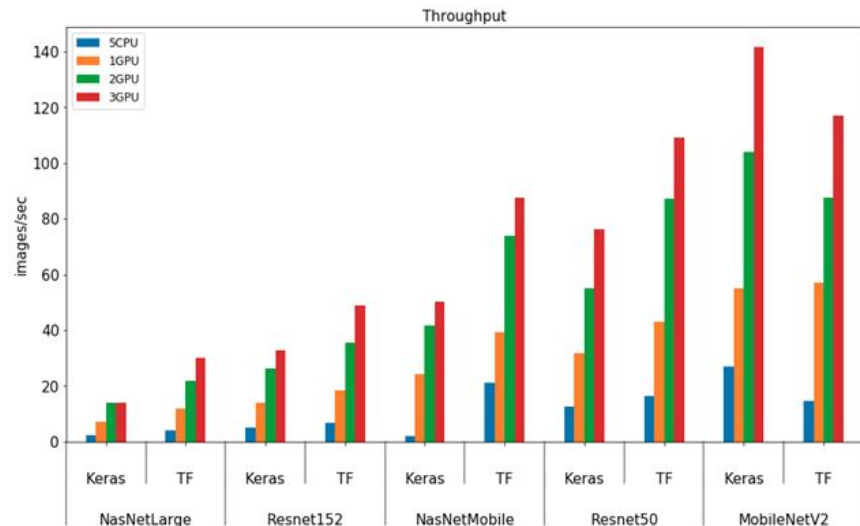
# AlexNet



**GPU1**

**GPU2**

- 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)
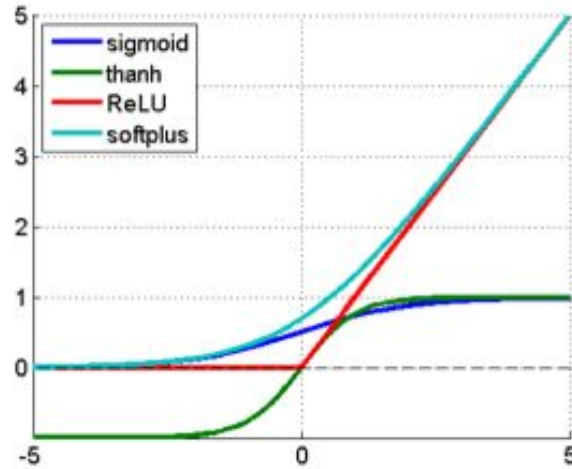    - Regularization (e.g. dropout)

    ....

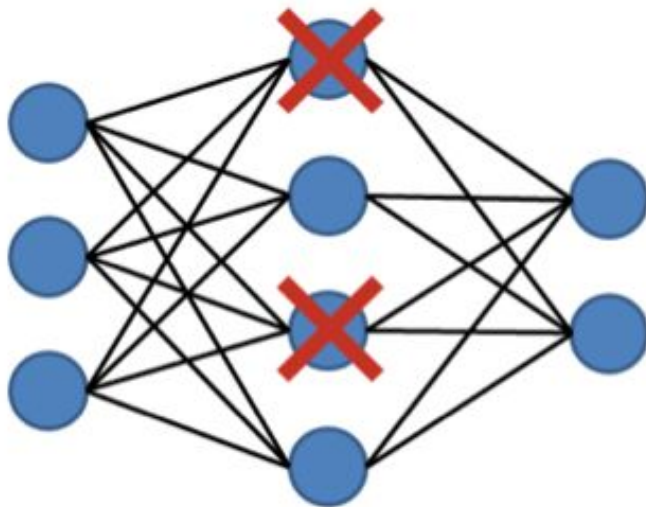# Rectified Linear Units  -  Activation function (2010)

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



- 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)

https://www.cs.toronto.edu/~fritz/absps/reluICML.pdf

# 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
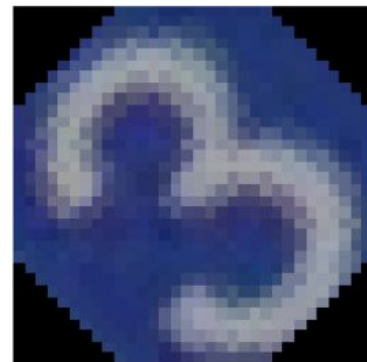
Original    Horizontal Flip    Pad & Crop    Rotate

# Training a neural network

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

- Predict gender from weight and height

# 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)
- **y** : variable being predicted (Gender)
- $y_{true}$ : true value of y, $y_{pred}$ : predicted value of y = **o**
- Loss function **L: MSE**
- Activation function **f**

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

- Outputs of the hidden layer **h**
- Unknown parameters: weights **w** and biases **b**

Input Layer          Hidden Layer          Output Layer

b1

W1

h1

W5          b3

W2

o1

W3          b2

gender

W6

weight

height

h2

W4

# Back propagation

- **Training the network == trying to minimize its loss**
  - Find weights **w** and biases **b**
  - $L(w_1, w_2, w_3, w_4, w_5, w_6, b_1, b_2, b_3)$
- Minimization taking partial derivatives  (**back propagation**)

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

For very simple case: with only Alice in the dataset, n=1

$$L = (1 - y_{pred})^2 \qquad y_{pred} = o_1 = f(w_5 h_1 + w_6 h_2 + b_3) \qquad h_1 = f(w_1 x_1 + w_2 x_2 + b_1)$$

$$\frac{\partial L}{\partial y_{pred}} = \frac{\partial (1 - y_{pred})^2}{\partial y_{pred}} \qquad \frac{\partial y_{pred}}{\partial h_1} = \boxed{w_5 * f'(w_5 h_1 + w_6 h_2 + b_3)} \qquad \frac{\partial h_1}{\partial w_1} = \boxed{x_1 * f'(w_1 x_1 + w_2 x_2 + b_1)}$$

27

# Gradient Descent

- optimization algorithm to find weights and biases to minimize loss

- **update equation**:

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



- **η (learning rate)** is a constant 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.
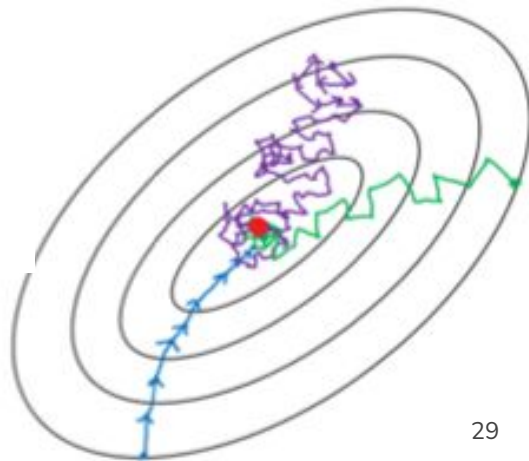
28

# Stochastic gradient descent (SDG)

- **Stochastic** -> the parameters are updated using only a single training instance (usually randomly selected) in each iteration

- Use mini-batch **sampled** in the dataset for gradient estimate.

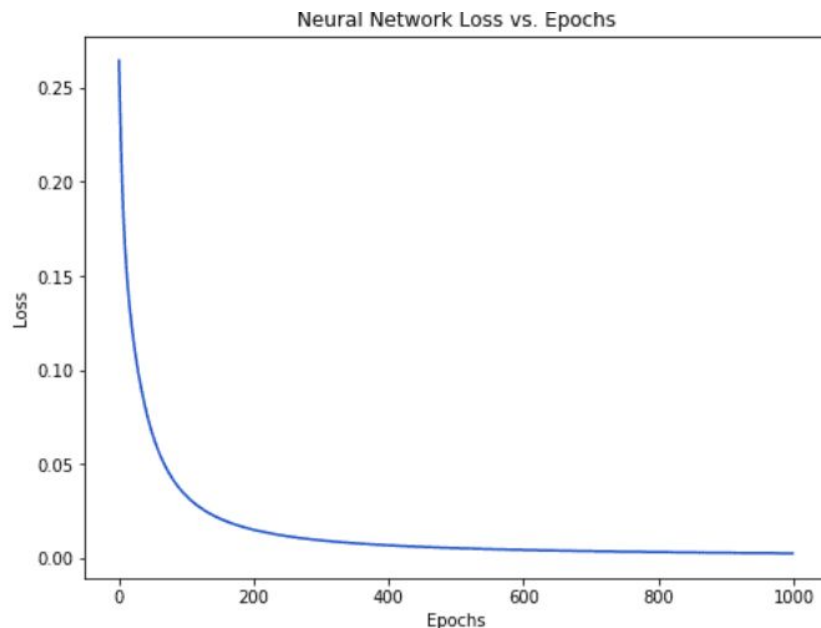$$\Theta^{t+1} = \Theta^t - \frac{\eta_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_\Theta \mathcal{L}_i$$

- 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

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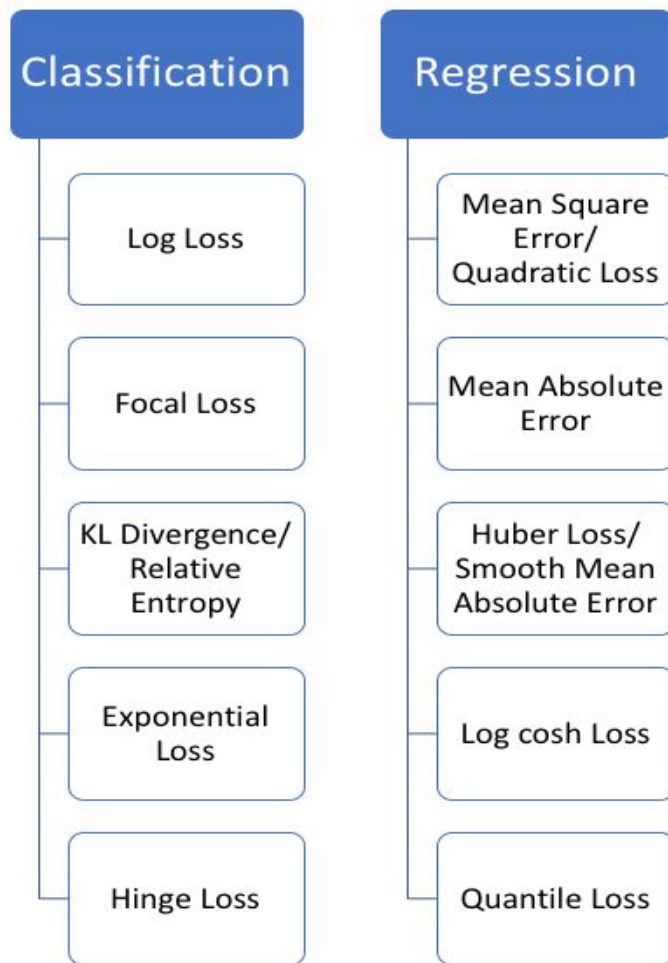— Stochastic gradient descent

# Training the network

- Choose one sample from our dataset
  - This is what makes it stochastic gradient descent - only operate on one sample at a time
- Calculate all the partial derivatives of loss with respect to weights or biases
- Use the update equation to update each weight and bias
- Iterate



Neural Network Loss vs. Epochs

# 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

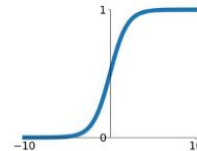| Classification | Regression |
| --- | --- |
| Log Loss | Mean Square Error/ Quadratic Loss |
| Focal Loss | Mean Absolute Error |
| KL Divergence/ Relative Entropy | Huber Loss/ Smooth Mean Absolute Error |
| Exponential Loss | Log cosh Loss |
| Hinge Loss | Quantile Loss |

# 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
    $$Min(\sum_{i=1}^{n}(y_i - w_i x_i)^2 + p\sum_{i=1}^{n}|w_i|)$$
  - **L2** or Ridge regularization adds a penalty which is the sum of the squared values of weights
    - L2+MSE
    $$Min(\sum_{i=1}^{n}(y_i - w_i x_i)^2 + p\sum_{i=1}^{n}(w_i)^2)$$
- **Dropout** in NN context: hidden nodes are dropped randomly
- **Early Stopping** is a time regularization technique which stops training based on given criteria

# 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
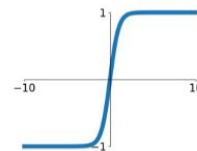- As a rule of thumb, start with ReLU
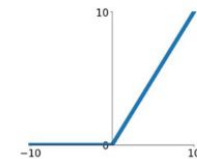
**Sigmoid**
$\sigma(x) = \frac{1}{1+e^{-x}}$
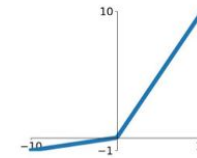
**tanh**
$\tanh(x)$

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