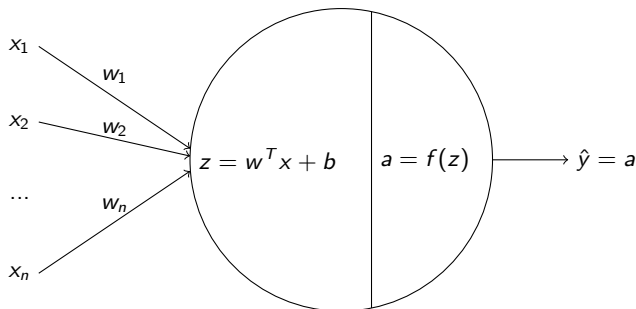# Neural networks
# Cours 1/3

**Machine Learning**
ModIA 2022

*Enseignant: Sixin Zhang*
`sixin.zhang@toulouse-inp.fr`

- *Pattern Recognition and Machine Learning*, Christopher M. Bishop - 2006

- *Deep Learning*, I Goodfellow, Y Bengio, A Courville - 2016

- *Understanding machine learning: From theory to algorithms*,
  S Shalev-Shwartz, S Ben-David - 2014

Ce cours a été conçu avec Sandrine Mouysset et Axel Carlier.

Representation of Perceptron



1. Inner product between input vector $x \in \mathbb{R}^n$ and the weight $w : w^T x$;
2. Add a bias scalar $(b \in \mathbb{R}) : z = w^T x + b$
3. Application of an activation function to $z : a = f(z)$
4. Output value $\hat{y} = a$, e.g. $\hat{y} \in \{0, 1\}$ for binary classification.

### Activation functions

The **activation functions**, denoted $f$, are usually non-linear functions. They can play a role of thresholding with 3 regimes,

- non-active: if the input value is under a threshold;
- transition phase: if the input value is close to the threshold;
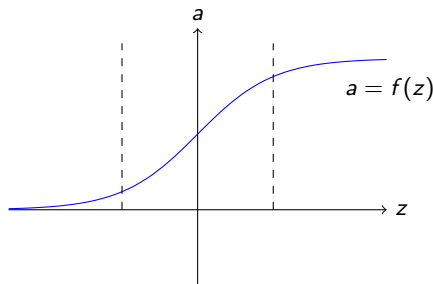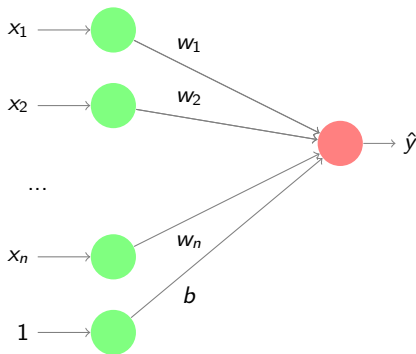- active: if the input value is above the threshold;



Figure: Sigmoid activation function : $f(z) = \frac{1}{1+e^{-z}}$

Convert input vector to $(x_1, \cdots, x_n, 1)^T \in \mathbb{R}^{n+1}$

$$w^T x + b = (w_1, \cdots, w_n, b)^T (x_1, \cdots, x_n, 1)$$

1. Initial weight $w^{(0)} = (w_i^{(0)})_{i \leq n}$
2. Draw training samples $(x^{\{1\}}, y^{\{1\}}), ..., (x^{\{m\}}, y^{\{m\}})$.
3. Compute the output of Perceptron and a loss $J(w)$:

$$\hat{y}^{\{j\}}(w) = f\left(\sum_{i=1}^{n} w_i x_i^{\{j\}}\right) \text{ and } J(w) = \frac{1}{m}\sum_{j=1}^{m} \ell(\hat{y}^{\{j\}}(w), y^{\{j\}})$$
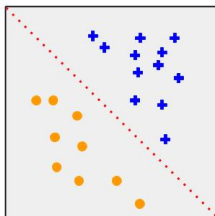
4. Update the weights from $w^{(t)}$ to $w^{(t+1)}$

$$w_i^{(t+1)} = w_i^{(t)} - \alpha^{(t)}\frac{\partial J}{\partial w_i}(w^{(t)})$$

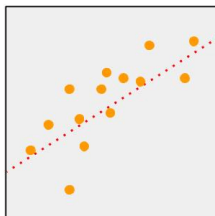   where $\alpha^{(t)}$ is a step size (learning rate) $\alpha^{(t)} > 0$.
5. Repeat 2-4 until convergence of $w^{(t)}$ or $J(w^{(t)})$.

$\Rightarrow$ How to define the **cost function** $\ell$ ?

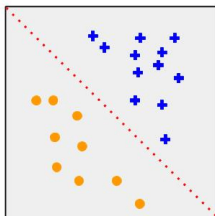Classification and regression



Classification

Regression

**Classification (Logistic regression)** Assign a category to each observation
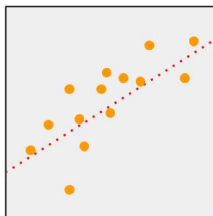**Binary case :** false/true, $y \in \{0, 1\}, \hat{y} \in [0, 1]$

- sigmoid activation ($\mathbb{R} \rightarrow [0, 1]$): $f(z) = (1 + e^{-z})^{-1}$
- Loss function: logistic cost (cross-entropy):

$$\text{loss}(\hat{y}, y) = -y log(\hat{y}) - (1 - y) log(1 - \hat{y})$$

Classification and regression



Classification                    Regression

**Linear Regression** Predict a real value of each observation :
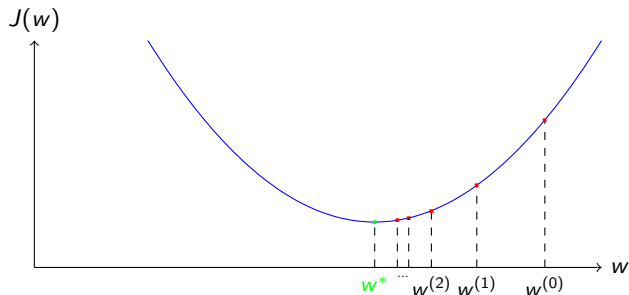
- linear activation : $f(z) = z$
- Mean squared error cost function (MSE):

$$\ell(\hat{y}, y) = (y - \hat{y})^2$$

$\Rightarrow$ How **to solve this type of problem** ?

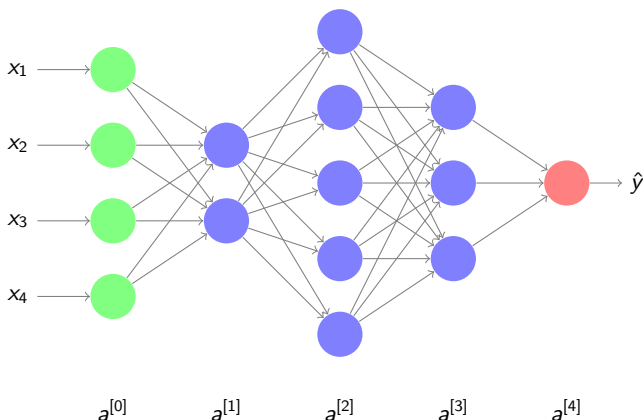**Gradient descent method: iterative method to find an optimal $w^*$**

1. Let $\hat{y}(w) = f\left(\sum_{i=1}^{n} w_i x_i\right)$ and $J(w) = \ell(\hat{y}(w), y)$
   - Assume $f$ is sigmoid.
   - What is the gradient of $J(w)$ with respect to $w_i$ in the following 2 cases?
     - $\ell(\hat{y}, y) = -y log(\hat{y}) - (1 - y) log(1 - \hat{y})$
     - $\ell(\hat{y}, y) = (y - \hat{y})^2$

2. Classification of a training set in 2d ($n = 2$) with $m = 4$ samples:

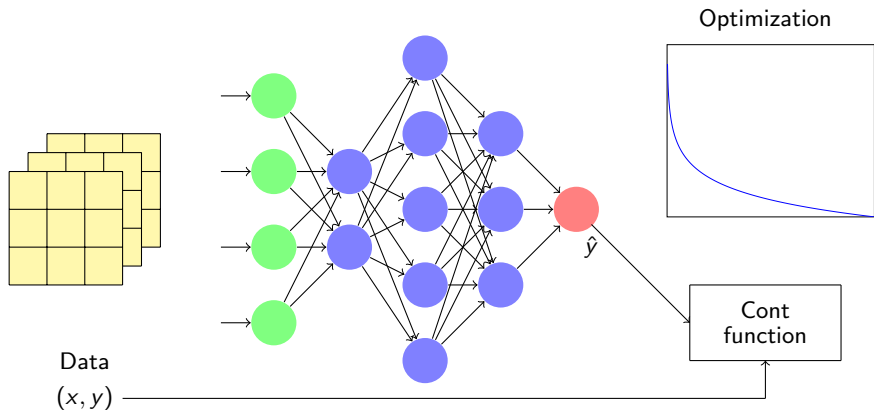   | $x_1$ | $x_2$ | $y$ |
   |-------|-------|-----|
   | 2     | 1     | 1   |
   | 0     | $-2$  | 1   |
   | $-2$  | 1     | 0   |
   | 0     | 2     | 0   |

   - Represent the data in 2d plane.
   - Propose a weight $w = (w_1, w_2)$ (without bias) such that the perceptron separates the training set with the sigmoid activation.
   - What is the minimal $J(w)$ with the logistic cost $\ell$?

A multi-layer perceptron (MLP) is composed of an input layer, several hidden layers and an output layer.

The **depth** of the network above is $L = 4$ (3 hidden layers plus one output layer).

Optimization
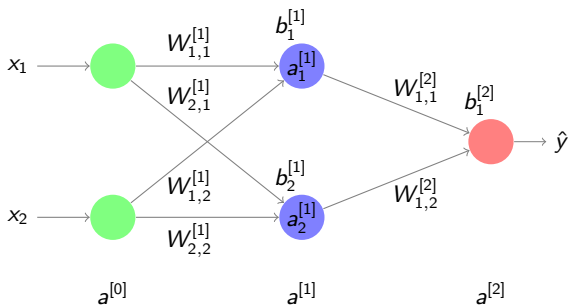
Data
$(x, y)$

$\hat{y}$

Cont
function

Multi-layer perceptron :

1. Functionality
2. Interpretation
3. Activation function
4. Multi-class classification loss

In order to train a multi-layer perceptron, we need to understand the following computational steps:

1. **Forward propagation** of input data to output;

2. Compute a **loss** from the output;

3. **Back propagation**: compute **gradients** of the loss with respect to the weights of the **output layer** and **hidden layers**;

4. **Update** all the weights based on optimization methods.

The weights of layer $k$: $W_{i,j}^{[k]}$ and $b_i^{[k]}$, $i$ output index, $j$ input index. For depth $L$, we denote all the weights by $\theta = (W^{[k]}, b^{[k]})_{k \leq L}$, e.g. $L = 2$

$$\hat{y}(x, \theta) = f \circ f^{[2]} \left( W^{[2]} f^{[1]}(W^{[1]} x + b^{[1]}) + b^{[2]} \right)$$

For an input $x^{\{i\}}$, we write the output $\hat{y}^{\{i\}}(\theta) = \hat{y}(x^{\{i\}}, \theta)$

2) Compute the objective function after the forward-propogation:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \ell(y^{\{i\}}, \hat{y}^{\{i\}}(\theta))$$
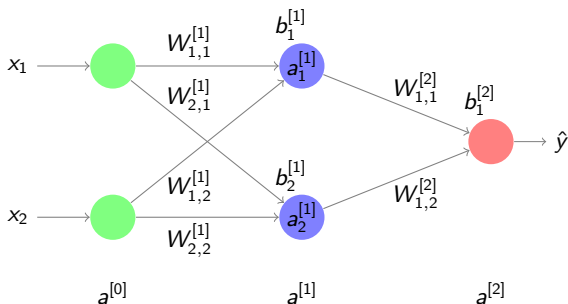
3) Back-propagation: to compute the gradients $\nabla_\theta J = (\frac{\partial J}{\partial \theta})^\mathsf{T}$ from output to input by the *chain rule* in Calculus, e.g.

$$\nabla_\theta J = \frac{1}{m} \sum_{i=1}^{m} \left( \frac{\partial \hat{y}^{\{i\}}}{\partial \theta} \right)^\mathsf{T} \nabla_{\hat{y}^{\{i\}}} \ell(y^{\{i\}}, \hat{y}^{\{i\}})$$

- Step 1: compute $\nabla_{\hat{y}^{\{i\}}} \ell(y^{\{i\}}, \hat{y}^{\{i\}})$ for $1 \le i \le m$.
- Step 2: compute $\frac{\partial \hat{y}^{\{i\}}}{\partial \theta}$ for $1 \le i \le m$.
- Step 3: compute $\nabla_\theta J$.

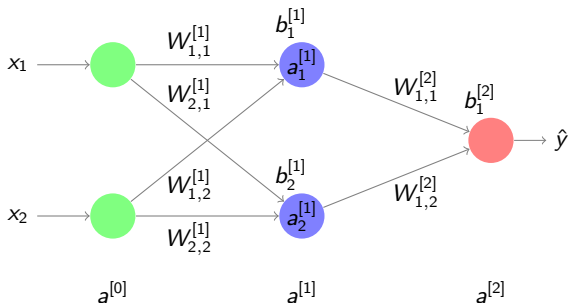Step 1 can be solved as in the previous Quiz. How about Step 2?

Assume $\hat{y} = f(a^{[2]}) \in \mathbb{R}$, $a^{[2]} = f^{[2]}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]}) \in \mathbb{R}$.

Compute $\frac{\partial \hat{y}}{\partial \theta}$:

$$\frac{\partial \hat{y}}{\partial b_1^{[2]}} = \frac{\partial \hat{y}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial b_1^{[2]}} = f'(a^{[2]})f^{[2]\prime}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]})$$
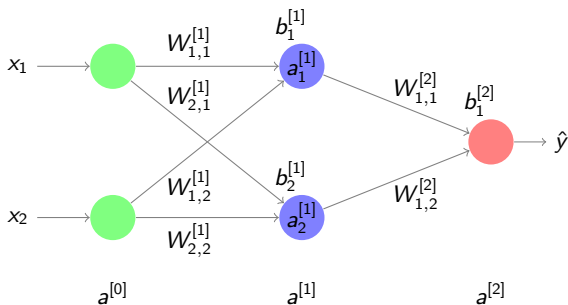
Assume $\hat{y} = f(a^{[2]}) \in \mathbb{R}$, $a^{[2]} = f^{[2]}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]}) \in \mathbb{R}$.
Compute $\frac{\partial \hat{y}}{\partial \theta}$:

$$\frac{\partial \hat{y}}{\partial W_{1,1}^{[2]}} = \frac{\partial \hat{y}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W_{1,1}^{[2]}} = f'(a^{[2]})f^{[2]'}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]})a_1^{[1]}$$
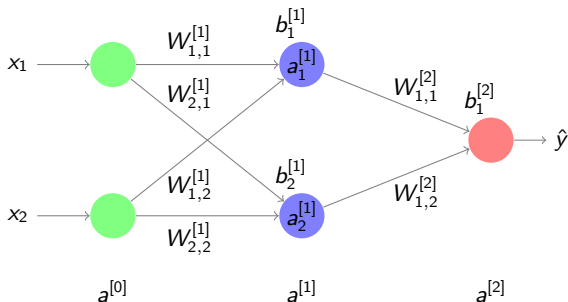
Assume $\hat{y} = f(a^{[2]}) \in \mathbb{R}$, $a^{[2]} = f^{[2]}(W_{1,1}^{[2]} a_1^{[1]} + W_{1,2}^{[2]} a_2^{[1]} + b_1^{[2]}) \in \mathbb{R}$.
Compute $\frac{\partial \hat{y}}{\partial \theta}$: Assume $a^{[1]} = f^{[1]}(W^{[1]} a^{[0]} + b^{[1]})$

$$\frac{\partial \hat{y}}{\partial W_{i,j}^{[1]}} = \frac{\partial \hat{y}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial W_{i,j}^{[1]}}$$

Jacobian matrices: $\frac{\partial a^{[1]}}{\partial W_{i,j}^{[1]}} : \mathbb{R}^1 \to \mathbb{R}^2$ $\frac{\partial a^{[2]}}{\partial a^{[1]}} : \mathbb{R}^2 \to \mathbb{R}^1$, $\frac{\partial \hat{y}}{\partial a^{[2]}} : \mathbb{R} \to \mathbb{R}$,

Assume $\hat{y} = f(a^{[2]}) \in \mathbb{R}$, $a^{[2]} = f^{[2]}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]}) \in \mathbb{R}$.

Compute $\frac{\partial \hat{y}}{\partial \theta}$: Assume $a^{[1]} = f^{[1]}(W^{[1]}a^{[0]} + b^{[1]})$

$$\frac{\partial a^{[2]}}{\partial a_1^{[1]}} = f^{[2]\prime}(W_{1,1}^{[2]}a_1^{[1]} + W_{1,2}^{[2]}a_2^{[1]} + b_1^{[2]})W_{1,1}^{[2]}$$

$$\frac{\partial a_i^{[1]}}{\partial W_{i,j}^{[1]}} = f_i^{[1]\prime}(W^{[1]}a^{[0]} + b^{[1]})a_j^{[0]}, \quad i,j \in \{1,2\}$$
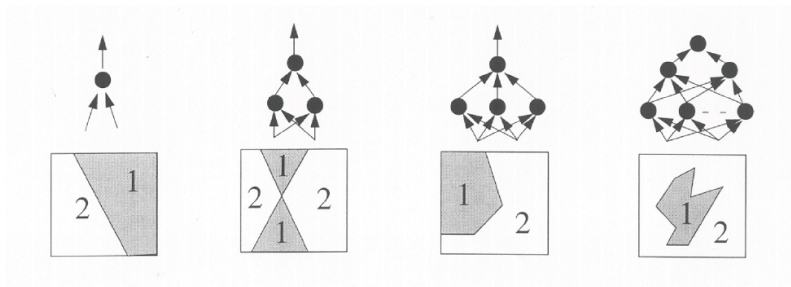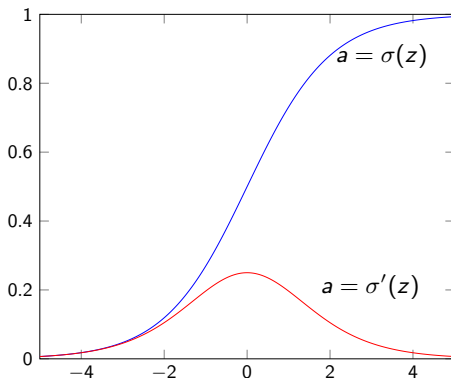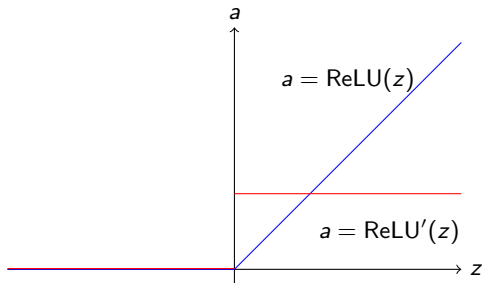
Figure: Linear vs. non-linear separation of training data

The non-linearity $f^{[1]}, f^{[2]}, \cdots$ in MLP plays a key role for non-linear separation.
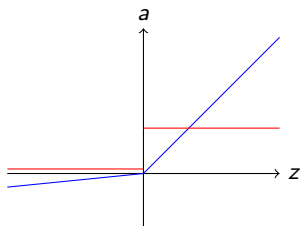
Example: https://playground.tensorflow.org/

- The gradient function tends to zero when $z$ is away from 0: cause vanishing gradients in the back propagation.
- Exercise: What is the relation between the sigmoid activation and the tanh activation: $\text{Tanh}(z) = (e^z - e^{-z})/(e^z + e^{-z})$.
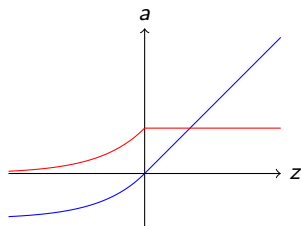
$$\text{ReLU}(z) = \left\{ \begin{array}{ll} z & \text{if} \quad z > 0 \\ 0 & \text{if} \quad z \leq 0 \end{array} \right.$$

The gradient is either 0 or 1 (when $z \neq 0$), very different to that of the sigmoid.

$$\text{leakyReLU}(z) = \begin{cases} \alpha z & \text{if } z < 0 \\ z & \text{if not} \end{cases} \qquad \text{eLU}(z) = \begin{cases} \alpha(e^z - 1) & \text{if } z < 0 \\ z & \text{if not} \end{cases}$$

*Leaky Rectified Linear Unit*        *Exponential Linear Unit*

There functions could potentially improve the gradient-descent training, e.g. to achieve faster convergence, but it is quite empirical.

- Question 1: Assume $y \in \{1, 2, \cdots, C\}$. To classify $x$ into $C$ categories, how to design a differentiable loss $\ell(y, \hat{y})$?
- Question 2: Assume $\hat{y}$ is a probability distribution over $\{1, 2, \cdots, C\}$, how to compute it from the output of an MLP?

- Answer 1: Use cross-entropy loss by representing $y$ and $\hat{y}$ as a probability distribution over $\{1, 2, \cdots, C\}$.
- Answer 2: Design a non-linear differentiable function $f$ such that $\hat{y} = f(a)$.

- KL divergence between two distributions $p$ and $q$ over $\{1, \cdots, C\}$

$$KL(q||p) = \sum_{i=1}^{C} \log \frac{q_i}{p_i} q_i$$

- Let $y$ be a vector in $\{0, 1\}^C$ such that $y_i = 1$ i.f.f the category of $y$ is $i$.
- Let $\hat{y}$ be a vector in $[0, 1]^C$ such that $\sum_i \hat{y}_i = 1$.
- The cross-entropy loss is $KL(y||\hat{y})$, which is equivalent to

$$\sum_{i=1}^{C} \log(y_i) y_i - \sum_{i=1}^{C} \log(\hat{y}_i) y_i$$

- In practice, we minimize the second term (to optimize MLP):

$$-\sum_{i=1}^{C} \log(\hat{y}_i) y_i$$

  - When $C = 2$, this loss is equivalent to the logistic regression loss.

- Let $a \in \mathbb{R}^C$ and $\hat{y} = f(a)$, such that

$$\hat{y}_i = \frac{e^{a_i}}{\sum_k e^{a_k}},$$

  - We have $\hat{y}_i \geq 0$ and $\sum_i \hat{y}_i = 1$.
- **Proposition**: for any $c \in \mathbb{R}$, $f(a + c) = f(a)$.
- To avoid numerical issues, we compute $f(a + c)$ with $c = -\max_j a_j$,

$$\hat{y}_i = \frac{e^{a_i - \max_j a_j}}{\sum_k e^{a_k - \max_j a_j}}.$$

- Exercise: Compute the derivative $\frac{\partial \hat{y}_i}{\partial a_k}$ for $i = k$ and $i \neq k$.
- Exercise (TP): MNIST classfication with MLP.