

Introduction to Deep Learning

Multi-Layered NNet and the back-propagation algorithm

Alexandre Allauzen



20/01/20

Outline

- 1 From logistic regression to neural network
- 2 From linear to non-linear classification
- 3 Multi-layered neural network and the back-propagation algorithm
- 4 Summary

Previously: logistic regression

Linear classification and the sigmoid function

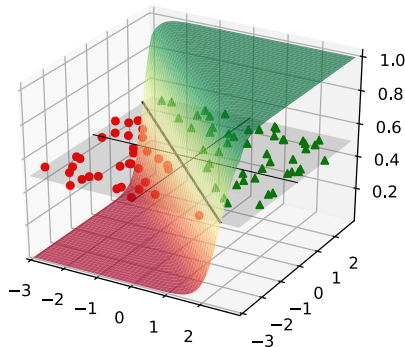
$$P(C = 1|\mathbf{x}) = \sigma(w_0 + \mathbf{w}^t \mathbf{x}) = y$$

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

$$a = w_0 + \mathbf{w}^t \mathbf{x}, a \in \mathbb{R}$$

$$a = w_0 + w_1 x_1 + \dots + w_D x_D$$

$$\boldsymbol{\theta} = (\mathbf{w}), \mathcal{D} = (\mathbf{x}_{(i)}, c_{(i)})_{i=1}^n$$



Learning with (Stochastic) Gradient Descent

$$\text{Minimize } \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})$$

$$l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)}) = \underbrace{-c_{(i)} \log(y_{(i)})}_{c_{(i)}=1} - \underbrace{(1 - c_{(i)}) \log(1 - y_{(i)})}_{c_{(i)}=0}$$

Bias or not bias

A matter of notation

The bias can be explicite:

$$\textcolor{red}{w}_0 + \mathbf{w}^t \mathbf{x} = w_0 + w_1 x_1 + \cdots + w_D x_D$$

or implicate:

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x} &= \begin{pmatrix} \textcolor{red}{w}_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{pmatrix} \cdot \begin{pmatrix} \textcolor{red}{1} \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \\ &= \textcolor{red}{w}_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 \end{aligned}$$

Gradient of the pointwise loss *w.r.t* \mathbf{w} - 1

$$\frac{\partial l(c_{(i)}, \mathbf{x}_{(i)}, \boldsymbol{\theta}_{(i)})}{\partial \mathbf{w}} = -c_{(i)} \frac{\partial \log y_{(i)}}{\partial \mathbf{w}} - (1 - c_{(i)}) \frac{\partial \log(1 - y_{(i)})}{\partial \mathbf{w}}$$

Setting $a = \mathbf{w}^t \mathbf{x}$, and consider the first term, the chain is

$$\mathbf{x} \rightarrow a \rightarrow y_{(i)} = \sigma(a) \rightarrow \log y_{(i)}$$

$$\begin{aligned} \frac{\partial \log y_{(i)}}{\partial \mathbf{w}} &= \frac{\partial \log y_{(i)}}{\partial y_{(i)}} \frac{\partial y_{(i)}}{\partial a} \frac{\partial a}{\partial \mathbf{w}} = \frac{1}{y_{(i)}} \frac{\partial y_{(i)}}{\partial a} \frac{\partial a}{\partial \mathbf{w}} \\ &= \frac{1}{y_{(i)}} \frac{\partial \sigma(a)}{\partial a} \frac{\partial a}{\partial \mathbf{w}} \\ &= \frac{1}{y_{(i)}} \sigma(a)(1 - \sigma(a)) \frac{\partial a}{\partial \mathbf{w}} = \frac{1}{y_{(i)}} y_{(i)}(1 - y_{(i)}) \frac{\partial a}{\partial \mathbf{w}} \\ &= \frac{1}{y_{(i)}} y_{(i)}(1 - y_{(i)}) \mathbf{x} = (1 - y_{(i)}) \mathbf{x} \end{aligned}$$

Gradient of the pointwise loss *w.r.t* \mathbf{w} - 2

Consider the second term

$$\begin{aligned}\frac{\partial \log(1 - y_{(i)})}{\partial \mathbf{w}} &= \frac{\partial \log(1 - y_{(i)})}{\partial y_{(i)}} \frac{\partial y_{(i)}}{\partial a} \frac{\partial a}{\partial \mathbf{w}} = \frac{-1}{1 - y_{(i)}} \frac{\partial y_{(i)}}{\partial a} \frac{\partial a}{\partial \mathbf{w}} \\ &= \frac{-1}{1 - y_{(i)}} y_{(i)} (1 - y_{(i)}) \mathbf{x} = -y_{(i)} \mathbf{x}\end{aligned}$$

Putting everything together:

$$\begin{aligned}\frac{\partial l(c_{(i)}, \mathbf{x}_{(i)}, \theta_{(i)})}{\partial \mathbf{w}} &= -c_{(i)}(1 - y_{(i)}) \mathbf{x} - (1 - c_{(i)})(-y_{(i)}) \mathbf{x} \\ &= -(c_{(i)} - c_{(i)}y_{(i)} - y_{(i)} + c_{(i)}y_{(i)}) \mathbf{x} \\ &= -(c_{(i)} - y_{(i)}) \mathbf{x}\end{aligned}$$

Update

Now consider the update, first for the class $c_{(i)} = 1$:

$$\begin{aligned}\mathbf{w}_{(i+1)} &= \mathbf{w}_{(i)} + \eta(c_{(i)} - y_{(i)})\mathbf{x} \\ &= \mathbf{w}_{(i)} + \eta(1 - y_{(i)})\mathbf{x}, \text{ note that } e_{(i)} = 1 - y_{(i)} \geq 0 \\ \mathbf{w}_{(i+1)}^t \mathbf{x} &= \mathbf{w}_{(i)}^t \mathbf{x} + \eta e_{(i)} \mathbf{x}^t \mathbf{x} > \mathbf{w}_{(i)}^t \mathbf{x} \\ \sigma(\mathbf{w}_{(i+1)}^t \mathbf{x}) &\geq \sigma(\mathbf{w}_{(i)}^t \mathbf{x}) \text{ therefore } P(C = 1|\mathbf{x}) \nearrow\end{aligned}$$

And for the class $c_{(i)} = 0$:

$$\begin{aligned}\mathbf{w}_{(i+1)} &= \mathbf{w}_{(i)} + \eta(c_{(i)} - y_{(i)})\mathbf{x} \\ &= \mathbf{w}_{(i)} - \eta y_{(i)}\mathbf{x}, \text{ note that } e_{(i)} = y_{(i)} \geq 0 \\ \mathbf{w}_{(i+1)}^t \mathbf{x} &= \mathbf{w}_{(i)}^t \mathbf{x} - \eta e_{(i)} \mathbf{x}^t \mathbf{x} < \mathbf{w}_{(i)}^t \mathbf{x} \\ \sigma(\mathbf{w}_{(i+1)}^t \mathbf{x}) &\leq \sigma(\mathbf{w}_{(i)}^t \mathbf{x}) \text{ therefore } P(C = 1|\mathbf{x}) \searrow\end{aligned}$$

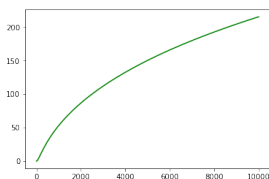
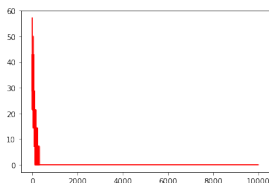
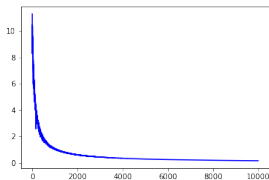
In both cases, the pointwise loss is reduced.

Unroll the updates

$$\begin{aligned}\mathbf{w}_{(1)} &= \mathbf{w}_{(0)} + \eta(c_{(1)} - y_{(1)})\mathbf{x}_{(1)} \\ \mathbf{w}_{(2)} &= \mathbf{w}_{(1)} + \eta(c_{(2)} - y_{(2)})\mathbf{x}_{(2)} \\ &= \mathbf{w}_{(0)} + \eta\left[(c_{(1)} - y_{(1)})\mathbf{x}_{(1)} + (c_{(2)} - y_{(2)})\mathbf{x}_{(2)}\right] \\ &\dots = \dots \\ \mathbf{w}_{(i)} &= \mathbf{w}_{(0)} + \eta \sum_{k=1}^i (c_{(k)} - y_{(k)})\mathbf{x}_{(k)}\end{aligned}$$

- \mathbf{w} is the sum of the training examples, weighted by the error.
- A kind of lossy compression of the \mathcal{D} .
- During training, the error tends to 0, so the update.
- What is the difference, between batch and online training ?

Training curves

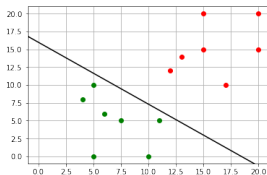


- Cumulated loss per epoch
- Classification error rate
- $\|\mathbf{w}\| + w_0^2$

The loss still decreases because the norm can increase (“forever”).

Omitted exercise

Assume we have w_0 and \mathbf{w} such that all the points are well classified.
Consider $w_0 = \alpha w_0$ and $\mathbf{w} = \alpha \mathbf{w}$, with $\alpha > 1$



- 1 what is the new classification error rate ?
- 2 does it change the decision boundary ?
- 3 what is the new loss function ?
- 4 propose a solution to overcome this issue, which is a kind of **overfitting**

l_2 regularisation

Consider the new loss function as a tradeoff, with $\lambda > 0$ and usually $\lambda < 1$ (like η):

$$\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2$$
$$\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathcal{D})}{\partial \boldsymbol{\theta}} = \nabla + \lambda \mathbf{w}$$

Therefore the update becomes

$$\begin{aligned} \mathbf{w}_{(i+1)} &= \mathbf{w}_{(i)} - \eta (\nabla + \mathbf{w}) \\ &= (1 - \eta \lambda) \mathbf{w}_{(i)} - \eta \nabla \end{aligned}$$

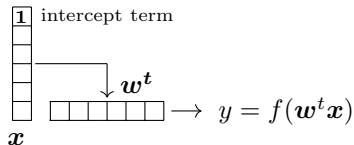
- A regularized objective function
- The new term is called regularization (l_2)
- Helps to avoid a kind of **overfitting**

Outline

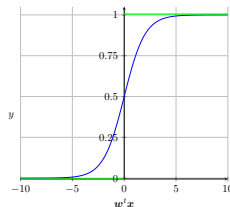
- 1 From logistic regression to neural network
- 2 From linear to non-linear classification
- 3 Multi-layered neural network and the back-propagation algorithm
- 4 Summary

A choice of terminology

Logistic regression (binary classification)

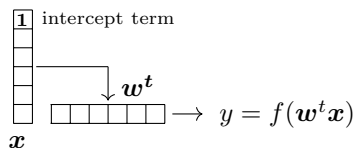


$$f(a = w^t x) = \frac{e^a}{1 + e^a} = \sigma(a)$$

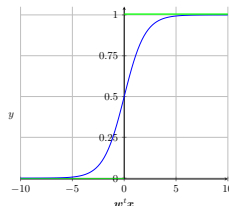


A choice of terminology

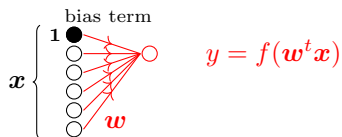
Logistic regression (binary classification)



$$f(a = \mathbf{w}^t \mathbf{x}) = \frac{e^a}{1 + e^a} = \sigma(a)$$



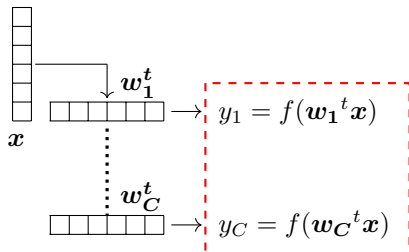
A single artificial neuron



- pre-activation : $a = \mathbf{w}^t \mathbf{x}$
- f : activation function
- Input values = input “neurones”
- \mathbf{x} : a vector of values, a layer

A choice of terminology - 2

From binary classification to C classes (Maxent)



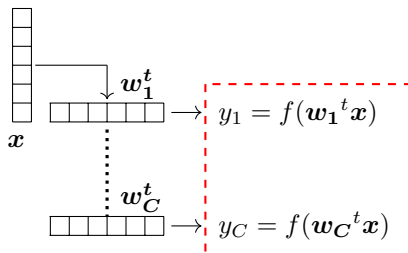
$$y_k = f(a_k = \mathbf{w}_k^t \mathbf{x}) = P(c = k | \mathbf{x})$$

$$= \frac{e^{a_k}}{\sum_{k'=1}^C e^{a_{k'}}} = \frac{e^{a_k}}{Z(\mathbf{x})}$$

$$\sum_{k=1}^C y_k = \sum_{k=1}^C P(c = k | \mathbf{x}) = 1$$

A choice of terminology - 2

From binary classification to C classes (Maxent)

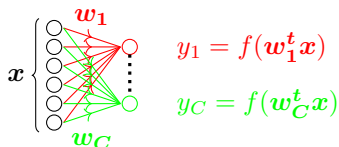


$$y_k = f(a_k = \mathbf{w}_k^t \mathbf{x}) = P(c = k | \mathbf{x})$$

$$= \frac{e^{a_k}}{\sum_{k'=1}^C e^{a_{k'}}} = \frac{e^{a_k}}{Z(\mathbf{x})}$$

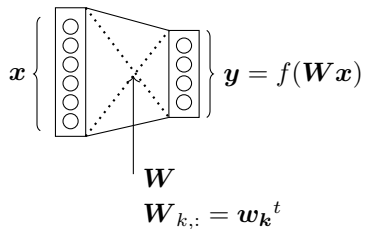
$$\sum_{k=1}^C y_k = \sum_{k=1}^C P(c = k | \mathbf{x}) = 1$$

A simple neural network

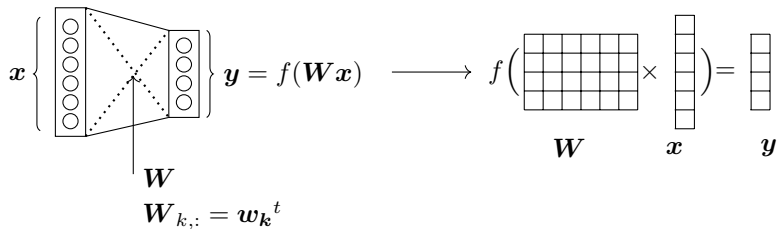


- x : input layer
- y : output layer
- each y_k has its parameters \mathbf{w}_k
- f is the **softmax** function

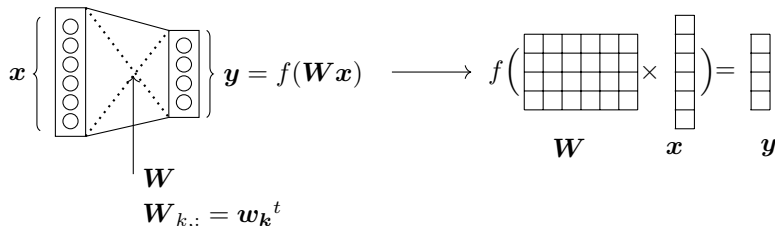
Two layers fully connected



Two layers fully connected



Two layers fully connected



Activation f :

- f is usually a non-linear function
 - f is a component wise function
 - tanh, sigmoid, relu, ...
- e.g* the softmax function:

Dimensions:

- $x : D \times 1$
- $W : C \times D$
- $y : (C \times \cancel{D}) \times (\cancel{D} \times 1) = C \times 1$

$$y_k = P(c = k | x) = \frac{e^{w_k^t x}}{\sum_{k'} e^{w_{k'}^t x}} = \frac{e^{W_{k,:} x}}{\sum_{k'} e^{W_{k',:} x}}$$

Matrix and Vector product

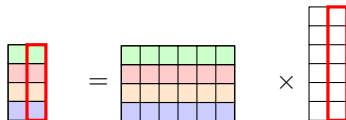
$$\begin{array}{ccccc}
 \mathbf{y} & = & \mathbf{W} & \times & \mathbf{x} \\
 \begin{array}{|c|} \hline \text{green} \\ \hline \text{red} \\ \hline \text{orange} \\ \hline \text{blue} \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|c|} \hline \text{green} & \text{green} & \text{green} & \text{green} & \text{green} \\ \hline \text{red} & \text{red} & \text{red} & \text{red} & \text{red} \\ \hline \text{orange} & \text{orange} & \text{orange} & \text{orange} & \text{orange} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \times & \begin{array}{|c|} \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \end{array} \\
 y_1 & = & \mathbf{W}_{1,:} & \times & \mathbf{x} \\
 y_2 & = & \mathbf{W}_{2,:} & \times & \mathbf{x} \\
 \dots & & \dots & & \dots
 \end{array}$$

In terms of dimension:

$$\text{With } \begin{cases} \mathbf{x} & : (L_1 \times 1) \\ \mathbf{W} & : (L_2 \times C_2) \\ \mathbf{y} & : (L_2 \times 1) \end{cases} \Rightarrow (\mathbf{W}\mathbf{x}) : (L_2 \times 1) = (L_2 \times \underbrace{\cancel{C_2}(\cancel{L_1})}_{C_2=L_1} \times 1)$$

Matrix-matrix product

\mathbf{X} is a matrix of 2 columns, 2 vectors as \mathbf{x} :

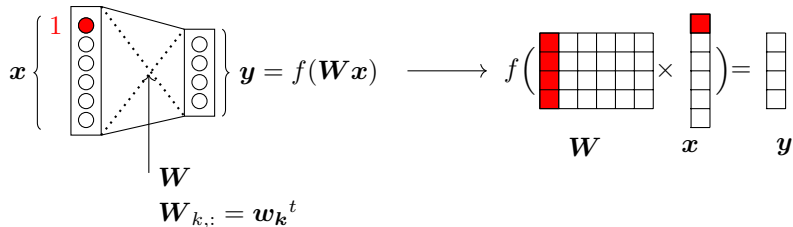
$$\mathbf{Y} = \mathbf{W} \times \mathbf{X}$$


In terms of dimension:

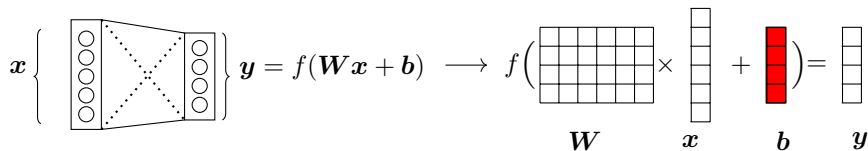
$$\text{With } \begin{cases} \mathbf{X} & : (L_1 \times C_1) \\ \mathbf{W} & : (L_2 \times C_2) \\ \mathbf{y} & : (L_2 \times C_1) \end{cases} \rightarrow (L_2 \times C_2) = (L_2 \times \underbrace{C_2}_{C_2=L_1}) \times (L_1 \times C_1)$$

Bias or not bias

Implicit Bias



Explicit bias



Classification with a simple neural network

Binary classification

- The input layer is a vector (\mathbf{x}), it encodes the data
- A single output neuron transform \mathbf{x} , \mathbf{w} are its parameters
- With a sigmoid activation, the loss function is the binary cross entropy, the log-loss, minus the log-likelihood, ...

Multiclass

- The input layer is a vector (\mathbf{x})
- The output layer \mathbf{y} contains one neuron per class
- It transforms \mathbf{x} , \mathbf{W} are the parameters of the output layer
- \mathbf{W} gathers the $\mathbf{w}_k = \mathbf{W}_{k,:}$
- With a softmax activation, the loss function is the cross entropy, the log-loss, minus the log-likelihood, ...

Other loss functions exist for classification

Regression with a simple neural network

Simple linear regression

$$y = f_{\theta}(\mathbf{x}), y \in \mathbb{R}$$

- The input layer is a vector (\mathbf{x}), it encodes the data
- A single output neuron for y (\mathbf{w} are its parameters)
- The activation function depends on the output domain

Multi-variate linear regression

$$\mathbf{y} = f_{\theta}(\mathbf{x}), \mathbf{y} \in \mathbb{R}^C$$

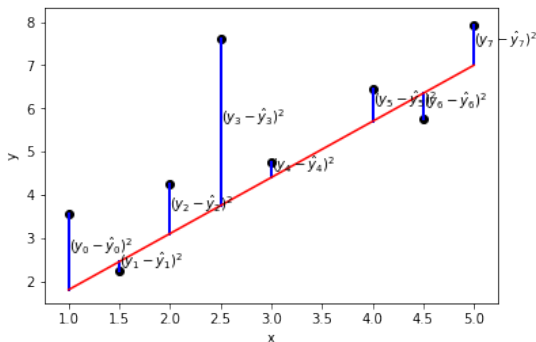
- The input layer is a vector (\mathbf{x})
- The output layer \mathbf{y} contains one neuron per output value
- It transforms \mathbf{x} in \mathbf{y} , \mathbf{W} are the parameters of the output layer

Loss for regression

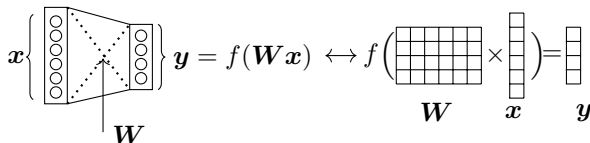
The Mean Squared Error *a.k.a* MSE

$$\mathcal{D} = (\mathbf{x}_{(i)}, \mathbf{y}_{(i)})_{i=1}^N,$$

$$\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_{(i)} - f_{\boldsymbol{\theta}}(\mathbf{x}_{(i)})\|^2$$



Two layers fully connected: another view

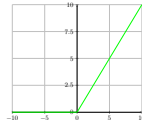
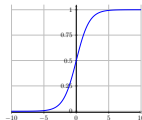
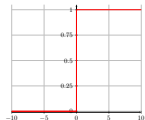


This basic brick implements a transformation of x in $y = f(Wx)$:

- A linear transformation Wx
- Followed by a non-linear function

Example: a candidate made 6 interviews $\rightarrow x \in \mathbb{R}^6$

- First compute 4 new scores : $Wx \in \mathbb{R}^4$, each is a linear combination of x
- Apply a non-linearity to get y



Outline

- 1 From logistic regression to neural network
- 2 From linear to non-linear classification
- 3 Multi-layered neural network and the back-propagation algorithm
- 4 Summary

Remember: one artificial neuron

Linear classification and the sigmoid function

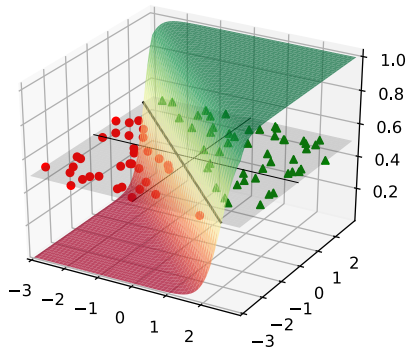
$$P(c = 1|\mathbf{x}) = \sigma(\mathbf{w}^t \mathbf{x}) = y$$

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

$$a = \mathbf{w}^t \mathbf{x}, a \in \mathbb{R}$$

$$a = w_0 + w_1 x_1 + \dots + w_D x_D$$

$$\boldsymbol{\theta} = (\mathbf{w}), \mathcal{D} = (\mathbf{x}_{(i)}, c_{(i)})_{i=1}^n$$

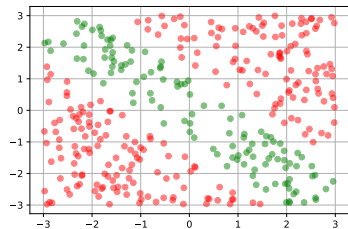
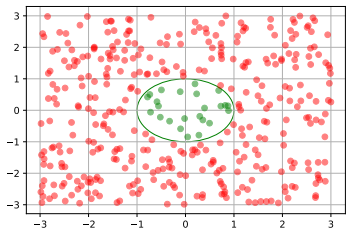
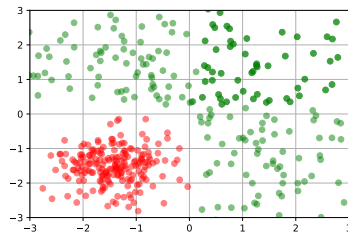
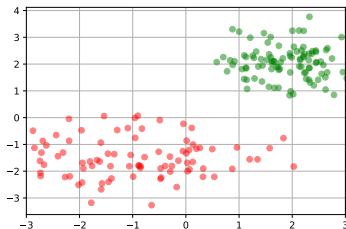


Learning with (Stochastic) Gradient Descent

$$\text{Minimize } \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})$$

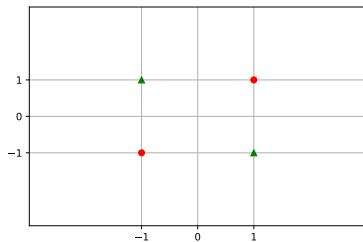
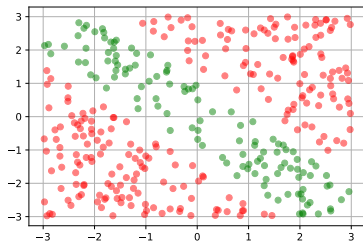
$$l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)}) = \underbrace{-c_{(i)} \log(y_{(i)})}_{c_{(i)}=1} - \underbrace{(1 - c_{(i)}) \log(1 - y_{(i)})}_{c_{(i)}=0}$$

Limits of the linear separation



Case study : X-or

Starting point

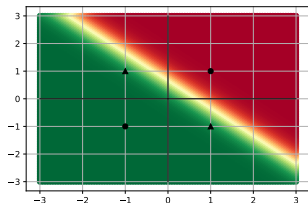


Hint

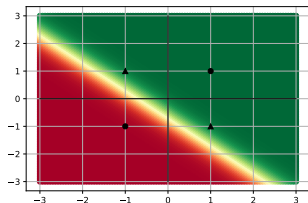
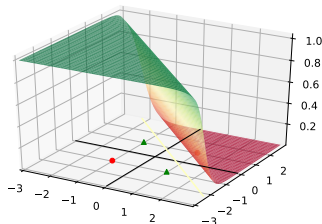
Try a first linear separation, maybe two.

Case study : X-or

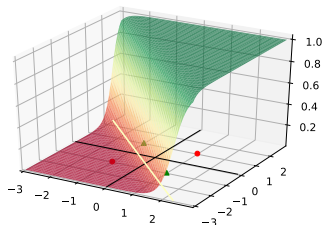
First try: linear classifiers



$$\mathbf{w} = [3, -4, -4]$$



$$\mathbf{w} = [2.5, 4, 4]$$



A first neural network

The multi-layer architecture

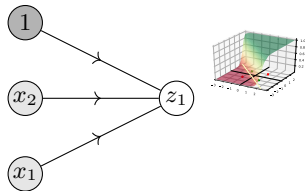
1

x_2

x_1

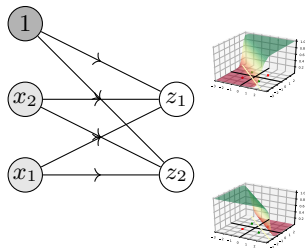
A first neural network

The multi-layer architecture



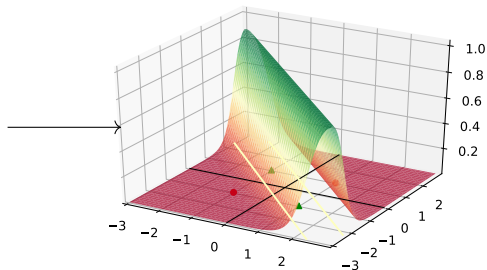
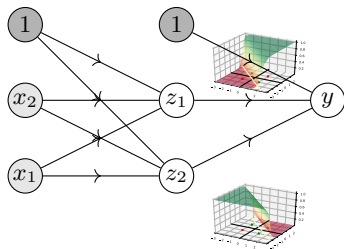
A first neural network

The multi-layer architecture



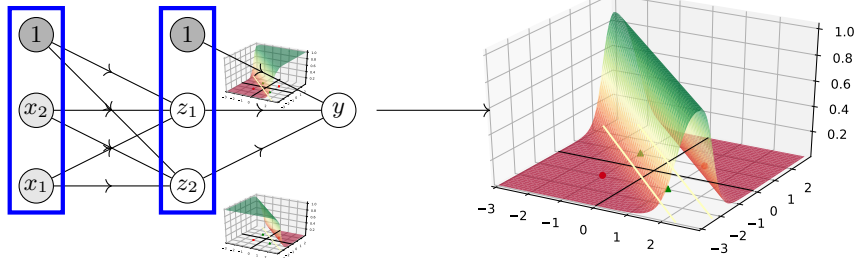
A first neural network

The multi-layer architecture



A first neural network

The multi-layer architecture



- The network is organized by layers: input (\mathbf{x}), hidden (\mathbf{z}), and output (\mathbf{y})
- Two layers are fully connected
- The propagation of the input is sequential:

$$z_1 = f(v_1^t \mathbf{x}) \text{ and } z_2 = f(v_2^t \mathbf{x})$$

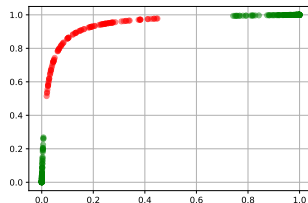
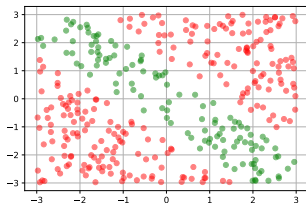
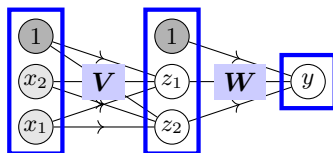
$$\Rightarrow \mathbf{z} = f(\mathbf{V}\mathbf{x})$$

$$y = f(\mathbf{w}^t \mathbf{z})$$

$$\Rightarrow y = f(\mathbf{W}\mathbf{z})$$

Another view of the neural network

Representation learning



- The network learns its own representation z .
- The final decision is linear.

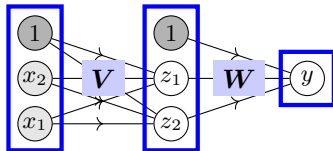
Summary

The multi-layer or feed-forward architecture

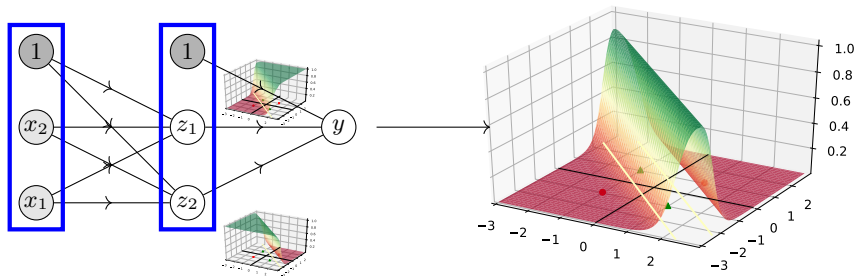
- 1 neuron = 1 value; 1 layer = 1 vector
- Two layers (\mathbf{x} , \mathbf{z}) fully connected:

$$\mathbf{z} = f(\mathbf{V}\mathbf{x})$$

- Inference: a sequential propagation
- Hidden layer (\mathbf{z}): the internal representation



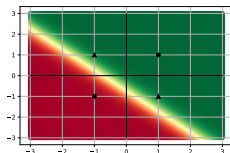
Exercise: the X-or network



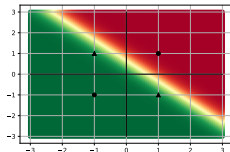
Compute the parameters of the output neuron (W) ?
 What is its activation function ?

Exercise : the X-or problem

From a boolean algebra point of view. Assume **true value is 1** and **false is 0**



- Write the truth table.
- Provide an interpretation as a boolean operator



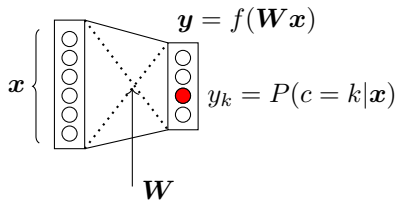
- Write the truth table.
- Provide an interpretation as a boolean operator

Write the boolean operation implemented by the previous neural network by merging the two previous steps.

Outline

- 1 From logistic regression to neural network
- 2 From linear to non-linear classification
- 3 Multi-layered neural network and the back-propagation algorithm
- 4 Summary

For a shallow network, with a single layer



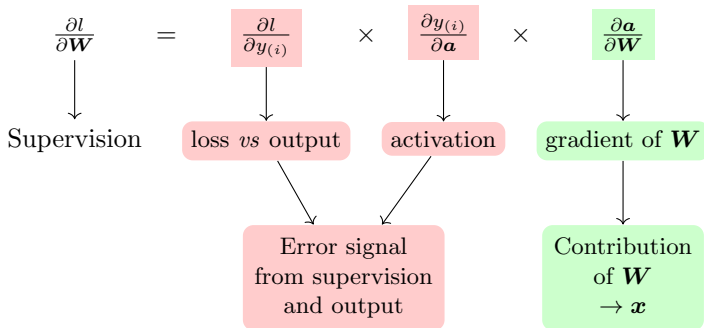
$$\theta = (W)$$

Forward (inference)	Backward (gradient)
$l(\theta, x_{(i)}, c_{(i)}) \leftarrow y$	$\frac{\partial l(\theta, x_{(i)}, c_{(i)})}{\partial W} = \frac{\partial l(\theta, x_{(i)}, c_{(i)})}{\partial y}$
$y = f(a)$	$\times \frac{\partial y}{\partial a}$
$a = Wx$	$\times \frac{\partial a}{\partial W}$

Review of the gradient computation

Without hidden layer (shallow net)

l is the shortcut for $l(c_{(i)}, \mathbf{x}_{(i)}, \boldsymbol{\theta}_{(i)})$

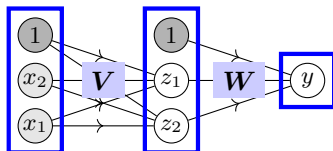


In the matrix form

$$\begin{aligned}
 \nabla_{\mathbf{W}} = \frac{\partial l}{\partial \mathbf{W}} &= \left(\frac{\partial l}{\partial y_{(i)}} \times \frac{\partial y_{(i)}}{\partial \mathbf{a}} \right) \times \left(\frac{\partial \mathbf{a}}{\partial \mathbf{W}} \right) \\
 &= \delta_{\mathbf{W}} \times \mathbf{x}^t \\
 &= (C, 1) \times (1, D) \\
 &= \text{Gradient up to the pre-activation} \times \text{input of the layer}
 \end{aligned}$$

- $\nabla_{\mathbf{W}}$ is a matrix of size (C, D)
- $\nabla_{\mathbf{W}}[i, j] = \delta[i] \times \mathbf{x}[j]$

Gradient computation with one hidden layer

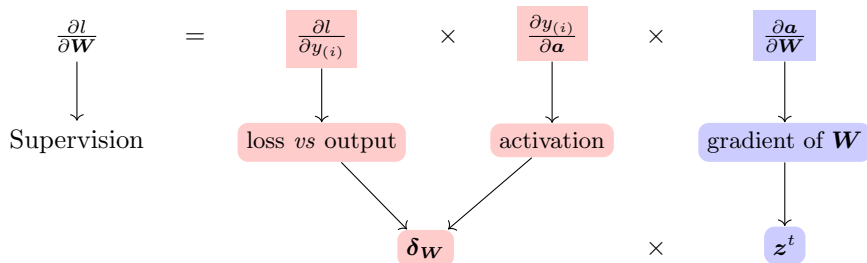


Two questions (or gradients):

- $\frac{\partial l}{\partial W} = ?$
- $\frac{\partial l}{\partial V} = ?$

Gradient computation with one hidden layer

Step 1: \mathbf{W}

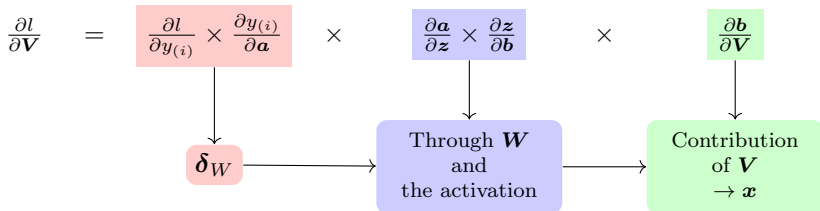


- Very similar to the shallow network, but the “input” is \mathbf{z} instead of \mathbf{x}
- But \mathbf{z} is computed by the network (parameters \mathbf{V})

Gradient computation with one hidden layer

Step 2 : \mathbf{V}

Denote the pre-activation of the input layer $\mathbf{b} = \mathbf{V}\mathbf{x}$



- The error signal δ_W is back-propagated through \mathbf{W} ,
- to get δ_V
- The update for \mathbf{V} is $\nabla_{\mathbf{V}} = \delta_V \mathbf{z}^t$

The back-propagation algorithm for a feed-forward network with one hidden layer

Introduced in (Rumelhart et al.1986)

One iteration of the online version, for a given value of θ :

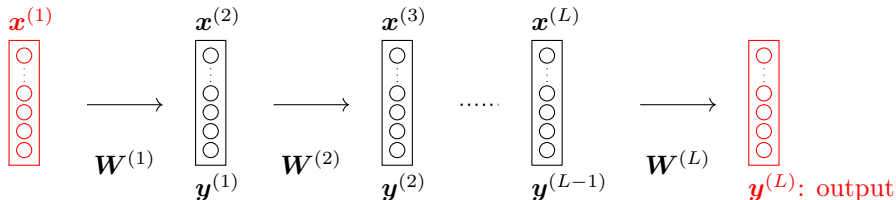
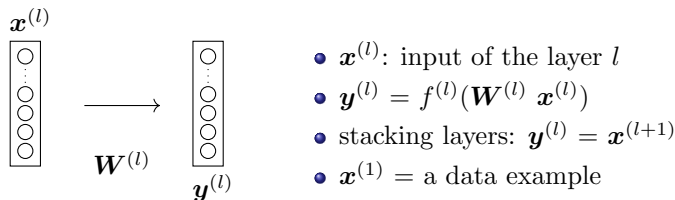
- ➊ Forward propagation of $\mathbf{x}_{(i)}$ ($\rightarrow \mathbf{z}$ and $\mathbf{y}_{(i)}$)
- ➋ Compute the loss
- ➌ Back-propagation and collect of the gradients:
 - output layer: δ_W and $\nabla_W = \delta_W \mathbf{z}^t$
 - input layer: δ_V and $\nabla_V = \delta_V \mathbf{x}^t$
- ➍ Update parameters:

$$\mathbf{W} = \mathbf{W} - \eta \nabla_W$$

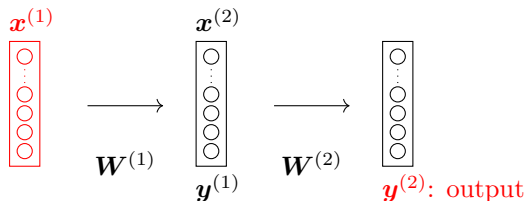
$$\mathbf{V} = \mathbf{V} - \eta \nabla_V$$

Notations for a multi-layer neural network (feed-forward)

One layer, indexed by l



Example : with one hidden layer



$$\theta = (W^{(1)}, W^{(2)})$$

To learn, we need the gradients for:

- the output layer: $\nabla_{W^{(2)}}$
- the hidden layer: $\nabla_{W^{(1)}}$

Back-propagation: generalization

For a hidden layer l :

- The gradient at the pre-activation level:

$$\delta^{(l)} = f'^{(l)}(\mathbf{a}^{(l)}) \circ (\mathbf{W}^{(l+1)^t} \delta^{(l+1)})$$

- The update is as follows:

$$\mathbf{W}^{(l)} = \mathbf{W}^{(l)} - \eta_t \delta^{(l)} \mathbf{x}^{(l)t}$$

The layer should keep:

- $\mathbf{W}^{(l)}$: the parameters
- $f^{(l)}$: its activation function
- $\mathbf{x}^{(l)}$: its input
- $\mathbf{a}^{(l)}$: its pre-activation associated to the input
- $\delta^{(l)}$: for the update and the back-propagation to the layer $l - 1$

Back-propagation: one training step

Pick a training example: $\mathbf{x}^{(1)} = \mathbf{x}_{(i)}$

Forward pass

For $l = 1$ to $(L - 1)$

- Compute $\mathbf{y}^{(l)} = f^{(l)}(\mathbf{W}^{(l)}\mathbf{x}^{(l)})$
- $\mathbf{x}^{(l+1)} = \mathbf{y}^{(l)}$

$$\mathbf{y}^{(L)} = f^{(L)}(\mathbf{W}^{(L)}\mathbf{x}^{(L)})$$

Backward pass

Init: $\delta^{(L)} = \nabla_{\mathbf{a}^{(L)}}$

For $l = L$ to 2 // all hidden units

- $\delta^{(l-1)} = f'^{(l-1)}(\mathbf{a}^{(l-1)}) \circ (\mathbf{W}^{(l)T} \delta^{(l)})$
- $\mathbf{W}^{(l)} = \mathbf{W}^{(l)} - \eta_t \delta^{(l)} \mathbf{x}^{(l)T}$

$$\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - \eta_t \delta^{(1)} \mathbf{x}^{(1)T}$$

Conclusion on back-propagation for one layer l

Training a NNet relies on forward-backward propagation.

Forward:

- get $\mathbf{x}^{(l)}$ for the previous layer;
- compute and send $\mathbf{y}^{(l)} = f^{(l)}(\mathbf{W}^{(l)}\mathbf{x}^{(l)})$.

Backward:

- get $\delta^{(l)}$ as input from the up-coming layer;
- compute and send $\delta^{(l-1)}$ to the previous layer;
- update parameters $\mathbf{W}^{(l)}$

Initialization recipes

A difficult question with several empirical answers.

One standard trick

$$\mathbf{W} \sim \mathcal{N}(0, \frac{1}{\sqrt{n_{in}}})$$

with n_{in} is the number of inputs

A more recent one

$$\mathbf{W} \sim \mathcal{U}\left[-\frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}, \frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}\right]$$

with n_{in} is the number of inputs

Outline

- 1 From logistic regression to neural network
- 2 From linear to non-linear classification
- 3 Multi-layered neural network and the back-propagation algorithm
- 4 Summary

Summary

Multi-layered Perceptron (MLP) or feed-forward NNet

- Artificial neurons are organized in **layers**: \rightarrow a vector
 - Two layers are in general **fully connected**
 - A linear transformation parametrized by a matrix
 - followed by a pointwise non-linear function, the activation function
 - A **feed-forward** architecture is a stack of layers fully connected
- \rightarrow Can approximate any functions depending on the number of hidden units.

Training by back-propagation

- After a forward pass (inference from input to the output)
- Backward pass (compute the gradients of each layer from the output to the input)



David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams.
1986.

Learning representations by back-propagating errors.

Nature, 323(6088):533–536, 10.