

# Neural Net: feed-forward architecture and back-propagation

A. Allauzen

Université Paris-Sud / LIMSI-CNRS



14 novembre 2017

# Outline

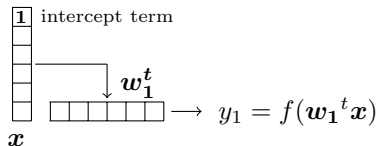
- 1 Neural Nets : Basics
  - Terminology
  - Training by back-propagation

# Outline

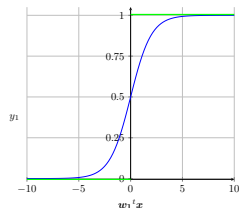
- 1 Neural Nets : Basics
  - Terminology
  - Training by back-propagation

# A choice of terminology

## Logistic regression (binary classification)

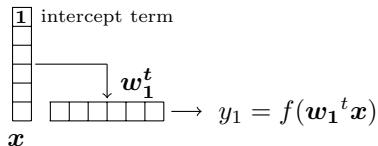


$$f(a = w_1^t \mathbf{x}) = \frac{1}{1 + e^{-a}}$$

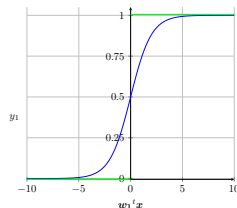


# A choice of terminology

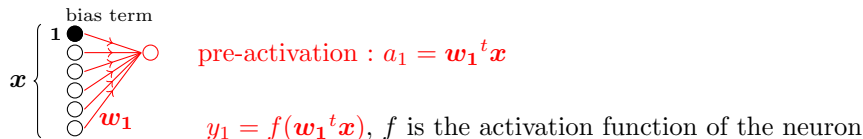
## Logistic regression (binary classification)



$$f(a = w_1^t x) = \frac{1}{1 + e^{-a}}$$

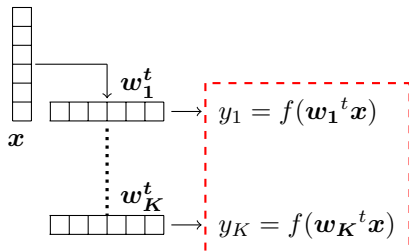


## A single artificial neuron



# A choice of terminology - 2

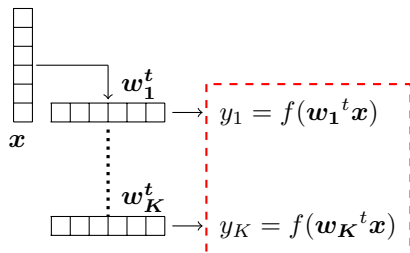
From binary classification to  $K$  classes (Maxent)



$$f(a_k = w_k^t x) = \frac{e^{a_k}}{\sum_{k'=1}^K e^{a_{k'}}} = \frac{e^{a_k}}{Z(x)}$$

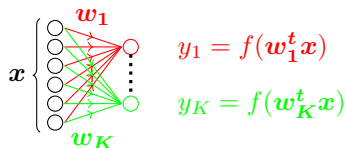
# A choice of terminology - 2

From binary classification to  $K$  classes (Maxent)



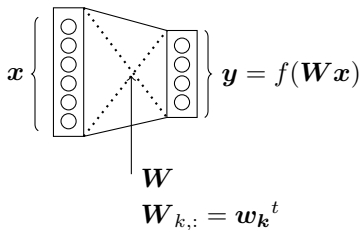
$$f(a_k = \mathbf{w}_k^t \mathbf{x}) = \frac{e^{a_k}}{\sum_{k'=1}^K e^{a_{k'}}} = \frac{e^{a_k}}{Z(\mathbf{x})}$$

A simple neural network



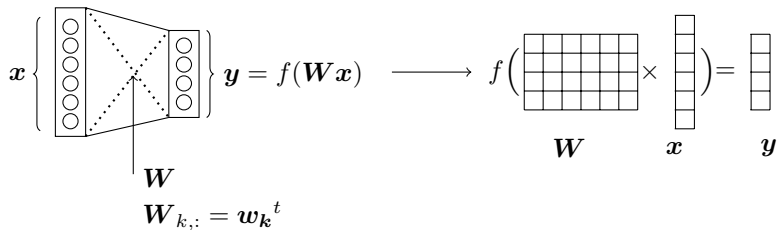
- $\mathbf{x}$  : input layer
- $\mathbf{y}$  : output layer
- each  $y_k$  has its parameters  $\mathbf{w}_k$
- $f$  is the **softmax** function

# Two layers fully connected

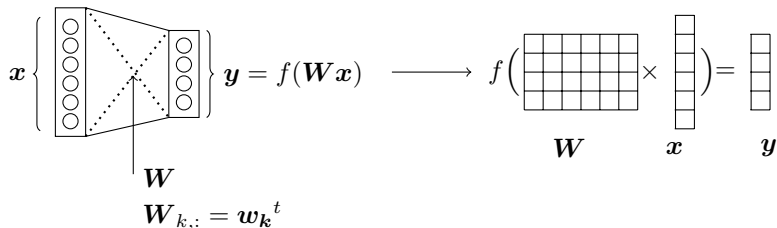




# Two layers fully connected



# Two layers fully connected



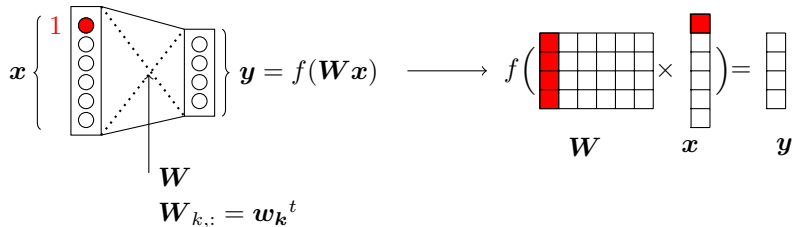
- $f$  is usually a non-linear function
- $f$  is a component wise function
- *e.g* the softmax function :

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

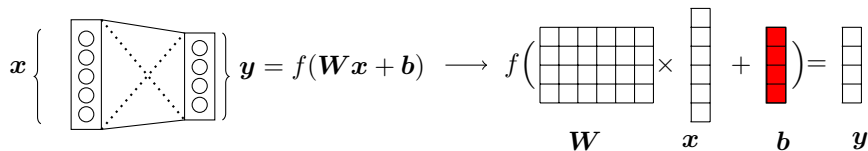
- tanh, sigmoid, relu, ...

# Bias or not bias

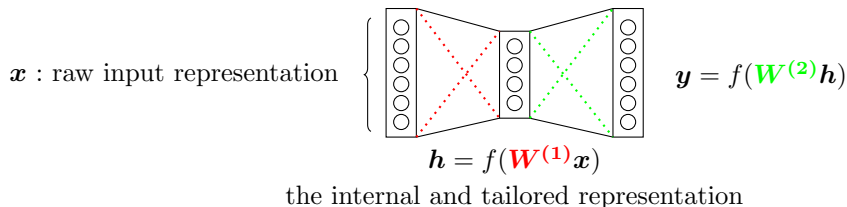
## Implicit Bias



## Explicit bias



# With neural network : add a hidden layer



## Intuitions

- Learn an internal representation of the raw input
- Apply a non-linear transformation
- The input representation  $\mathbf{x}$  is transformed/compressed in a new representation  $\mathbf{h}$
- Adding more layers to obtain a more and more abstract representation

# How do we learn the parameters ?

For a supervised single layer neural net

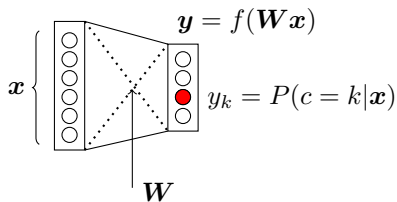
Just like a maxent model :

- Calculate the gradient of the objective function and use it to iteratively update the parameters.
- Conjugate gradient, L-BFGS, ...
- In practice : **Stochastic gradient descent (SGD)**

With one hidden layer

- The internal (“hidden”) units make the function non-convex ... just like other models with hidden variables :
  - hidden CRFs (?), ...
- But we can use the same ideas and techniques
- Just without guarantees  $\Rightarrow$  **backpropagation** (?)

# Ex. 1 : A single layer network for classification



$\theta$  = the set of parameters, in this case :

$$\theta = (\mathbf{W})$$

## The log-loss (conditional log-likelihood)

Assume the dataset  $\mathcal{D} = (\mathbf{x}_{(i)}, c_{(i)})_{i=1}^N$ ,  $c_{(i)} \in \{1, 2, \dots, C\}$

$$\mathcal{L}(\theta) = \sum_{i=1}^N l(\theta, \mathbf{x}_{(i)}, c_{(i)}) = \sum_{i=1}^N \left( - \sum_{c=1}^C \mathbb{I}\{c = c_{(i)}\} \log(P(c|\mathbf{x}_{(i)})) \right) \quad (1)$$

$$l(\theta, \mathbf{x}_{(i)}, c_{(i)}) = - \sum_{k=1}^C \mathbb{I}\{k = c_{(i)}\} \log(y_k) \quad (2)$$

## Ex. 1 : optimization method

### Stochastic Gradient Descent (?)

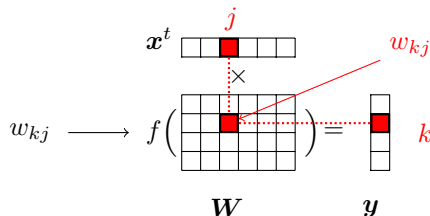
For (  $t = 1$  ; until convergence ;  $t++$  ) :

- Pick randomly a sample  $(\mathbf{x}_{(i)}, c_{(i)})$
- Compute the gradient of the loss function w.r.t the parameters  $(\nabla_{\boldsymbol{\theta}})$
- Update the parameters :  $\boldsymbol{\theta} = \boldsymbol{\theta} - \eta_t \nabla_{\boldsymbol{\theta}}$

### Questions

- convergence : what does it mean ?
- what do you mean by  $\eta_t$  ?
  - convergence if  $\sum_t \eta_t = \infty$  and  $\sum_t \eta_t^2 < \infty$
  - $\eta_t \propto t^{-1}$
  - and lot of variants like Adagrad (?), Down scheduling, ... see (?)

# Ex. 1 : compute the gradient - 1



Inference chain :

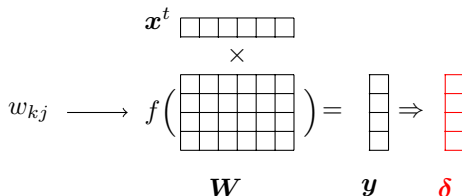
$$\mathbf{x}_{(i)} \longrightarrow (\mathbf{a} = \mathbf{W} \mathbf{x}_{(i)}) \longrightarrow (\mathbf{y} = f(\mathbf{a})) \longrightarrow l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})$$

The gradient for  $w_{kj}$

$$\begin{aligned} \nabla_{w_{kj}} &= \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial w_{kj}} = \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial \mathbf{y}} \times \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \times \frac{\partial \mathbf{a}}{\partial w_{kj}} \\ &= -(\mathbb{I}\{k = c_{(i)}\} - y_k) x_j = \delta_k x_j \end{aligned}$$



## Ex. 1 : compute the gradient - 2

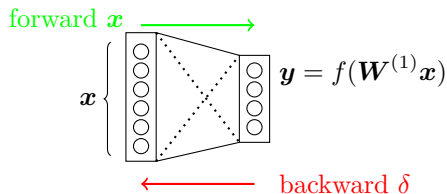


### Generalization

$$\begin{aligned}\nabla_{\mathbf{W}} &= \delta \mathbf{x}^t \\ \delta &= \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial \mathbf{a}} \\ \delta_k &= -(\mathbb{I}\{k = c_{(i)}\} - y_k)\end{aligned}$$

with  $\delta$  the gradient at the pre-activation level.

## Ex. 1 : Summary



### Inference : a forward step

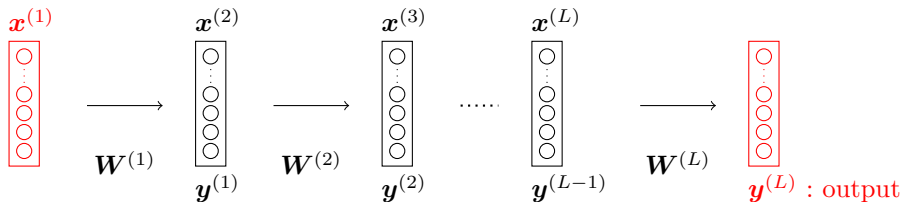
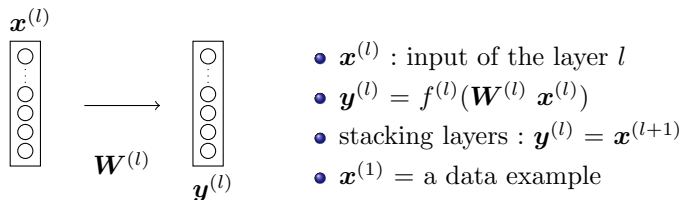
- matrix multiplication with the input  $x$
- Application of the activation function

### One training step : forward and backward steps

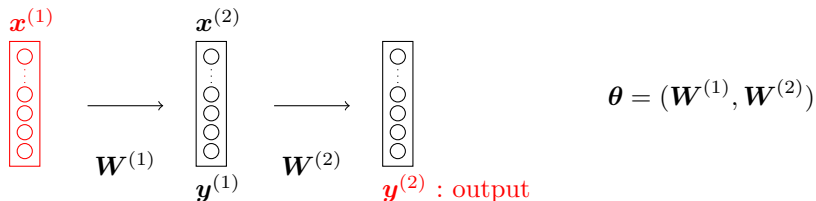
- Pick randomly a sample  $(x_{(i)}, c_{(i)})$
- Compute  $\delta$
- Update the parameters :  $\theta = \theta - \eta_t \delta x^t$

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

One layer, indexed by  $l$



## Ex. 2 : with one hidden layer



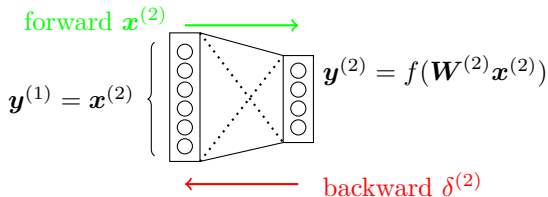
To learn, we need the gradients for :

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

# Back-propagation of the loss gradient

For the output layer

As in the Ex. 1 :



$$\nabla_{\mathbf{W}^{(2)}} = \delta^{(2)} \mathbf{x}^{(2)t}, \text{ with}$$

$$\delta_k^{(2)} = -(\mathbb{I}\{k = c_{(i)}\} - y_k)$$

$$\mathbf{y} \rightarrow \mathbf{y}^{(2)}$$

$$\mathbf{W} \rightarrow \mathbf{W}^{(2)}$$

$$\mathbf{x} \rightarrow \mathbf{x}^{(2)} = \mathbf{y}^{(1)}$$

# Back-propagation of the loss gradient

For the hidden layer - 1

The goal : compute  $\delta^{(1)}$

Inference (/forward) chain from  $\mathbf{a}^{(1)}$  to the output :

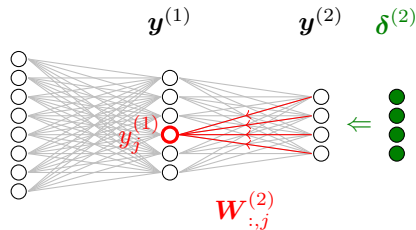
$$\mathbf{y}^{(1)} = f^{(1)}(\mathbf{a}^{(1)}) \rightarrow \left( \mathbf{a}^{(2)} = \mathbf{W}^{(2)} \mathbf{y}^{(1)} \right) \rightarrow \left( \mathbf{y}^{(2)} = f^{(2)}(\mathbf{a}^{(2)}) \right) \rightarrow l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})$$

Backward / Back-propagation :

$$\delta_j^{(1)} = \nabla_{a_j^{(1)}} = \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial a_j^{(1)}} = \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial \mathbf{y}^{(2)}} \times \frac{\partial \mathbf{y}^{(2)}}{\partial \mathbf{a}^{(2)}} \times \frac{\partial \mathbf{a}^{(2)}}{\partial y_j^{(1)}} \times \frac{\partial y_j^{(1)}}{\partial a_j^{(1)}}$$

# Back-propagation of the loss gradient

For the hidden layer - 2

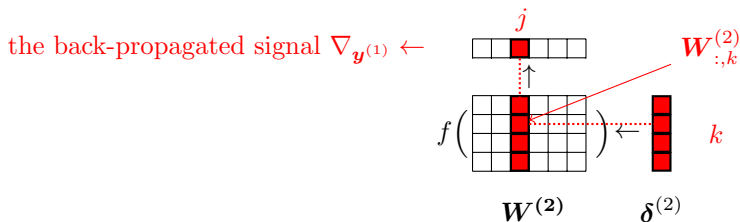


Backward / Back-propagation :

$$\begin{aligned}\delta_j^{(1)} = \nabla_{a_j^{(1)}} &= \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial a_j^{(1)}} = \frac{\partial l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)})}{\partial \mathbf{y}^{(2)}} \times \frac{\partial \mathbf{y}^{(2)}}{\partial \mathbf{a}^{(2)}} \times \frac{\partial \mathbf{a}^{(2)}}{\partial y_j^{(1)}} \times \frac{\partial y_j^{(1)}}{\partial a_j^{(1)}} \\ &= f'^{(1)}(a_j) \left( \mathbf{W}_{:,j}^{(2)T} \boldsymbol{\delta}^{(2)} \right)\end{aligned}$$

# Back-propagation of the loss gradient

For the hidden layer - 3



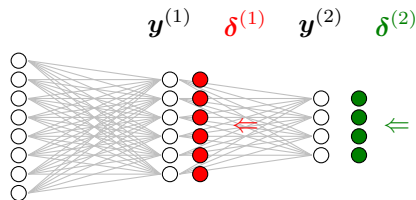
$$\nabla_{\mathbf{y}^{(1)}} = \mathbf{W}^{(2)t} \delta^{(2)}, \text{ then}$$

$$\delta^{(1)} = \nabla_{\mathbf{a}^{(1)}} = f^{(1)'}(\mathbf{a}^{(1)}) \circ (\mathbf{W}^{(2)t} \delta^{(2)})$$



# Back-propagation of the loss gradient

For the hidden layer - 4



As for the output layer, the gradient is :

$$\nabla_{\mathbf{W}^{(1)}} = \delta^{(1)} \mathbf{x}^{(1)t}, \text{ with}$$

$$\delta_j^{(1)} = \nabla_{a_j^{(1)}}$$

$$\delta^{(1)} = f'^{(1)}(\mathbf{a}^{(1)}) \circ (\mathbf{W}^{(2)t} \delta^{(2)})$$

The term  $(\mathbf{W}^{(2)t} \delta^{(2)})$  comes from the upper layer.

# 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