# Feed-Forward Neural Network

**Lecture 5**
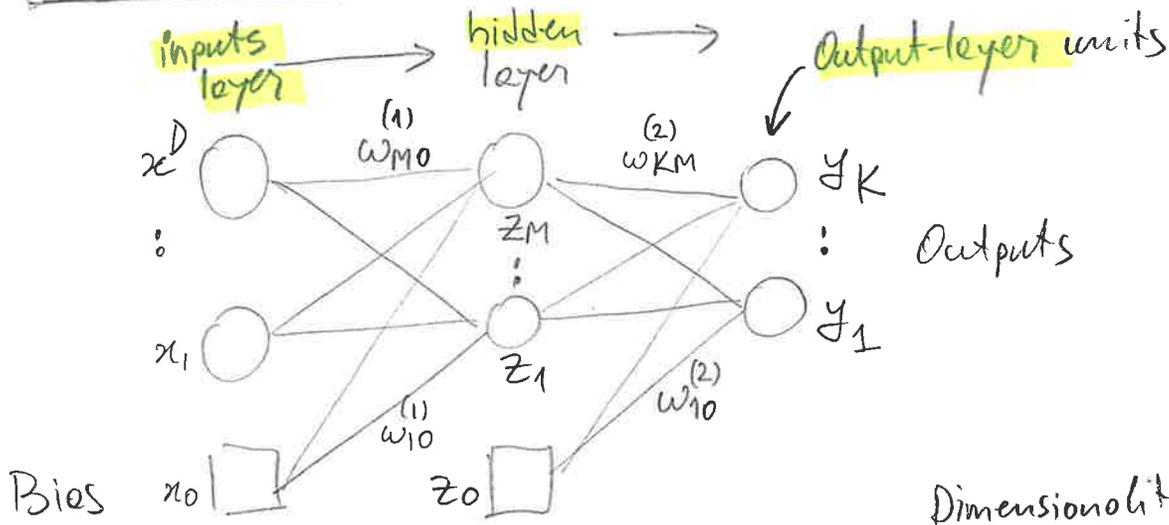


inputs layer → hidden layer → Output-layer units

Outputs

Bies

**Fig. 1**

Number units on hidden Layer

Dimensionality of input $x$

$$y_k(X, \omega) = \nabla \left( \sum_{m=0}^{M} \omega_{km}^{(2)} \cdot h\left( \sum_{d=0}^{D} \omega_{md}^{(1)} x_d \right) \right) \qquad (1)$$

Activation

$$h(a_m) = Z_m^{(1)} : \text{output of hidden units on Layer 1}$$

E.g.: Sigmoid function when Logistic Regression for Binary Classification

Denotes the non-linear activation function for a hidden unit.

≙ A **Neural Network** is a non-linear function that transforms the input $\boxed{x}$ into an output $\boxed{y}$ that is controlled by the set of parameters $\boxed{\omega}$.

==Training==

— Need objective/<u>loss</u>/<u>cost</u> function.

— For Linear Regression, can use |least squares loss|:

(2)
$$\mathcal{L}(\omega) = \frac{1}{N} \sum_{n=1}^{N} [y(x_n, \omega) - y_n]^2$$

— For Binary classification (uses <u>sigmoid</u> output activation function) the |negated log-likelihood| (or cross-entropy) is a typical loss function:

(3)
$$\mathcal{L}(\omega) = -\sum_{n=1}^{N} [y_n \cdot \log(\hat{y}_n) + (1-y_n) \log(1-\hat{y}_n)]$$

— For Multi-class classification problem (produced by <u>softmax</u> activ. function), we can use |negated-log-likelihood| (cross-entropy) loss function:

(4)
$$\mathcal{L}(\omega) = -\sum_{n=1}^{N} \sum_{k=1}^{K} y_{kn} \cdot \log\left[ \frac{e^{a_k(x,\omega)}}{\sum_{j=1}^{K} e^{a_j(x,\omega)}} \right]$$

**Backpropagation** = A procedure by which we pass errors backwards through a feed-forward NN in order to compute gradients for the weight parameters of the network.
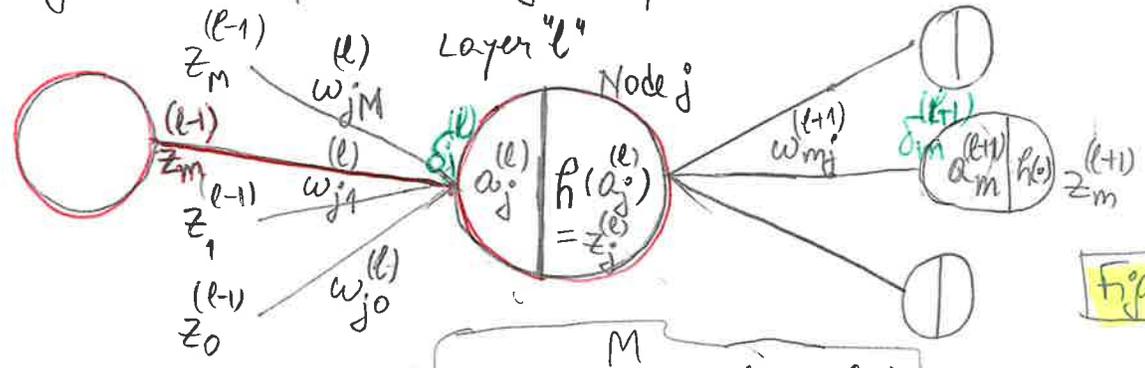


Fig. 2

**Activation:**
$$a_j^{(l)} = \sum_{m=0}^{M} w_{jm}^{(l)} \cdot z_m^{(l-1)} \quad (5)$$

— Output value: obtained by **activation function** $h(\cdot)$:

$$z_j^{(l)} = h(a_j^{(l)}) \quad (6)$$

— Computing derivatives of the objective function with respect to weights; particularly w.r.t. individual weight $w_{jm}^{(l)}$ ( $m^{th}$ weight for activation $j$ in layer $l$ ):

Use the chain rule:

$$(7) \quad \frac{\partial L}{\partial w_{jm}^{(l)}} = \underbrace{\frac{\partial L}{\partial a_j^{(l)}}}_{= \delta_j^{(l)}} \cdot \underbrace{\frac{\partial a_j^{(l)}}{\partial w_{jm}^{(l)}}}_{= z_m^{(l-1)} \text{ from eq. (5)}}$$

introduced notation for "errors"

$$(8) \quad \frac{\partial L}{\partial w_{jm}^{(l)}} = \delta_j^{(l)} \cdot z_m^{(l-1)}$$

← Derivative of Loss w.r.t. an arbitrary weight in network can be calculated as product of the error at the "output end of that weight" and value $z_m^{(l-1)}$ at "input end of the weight"

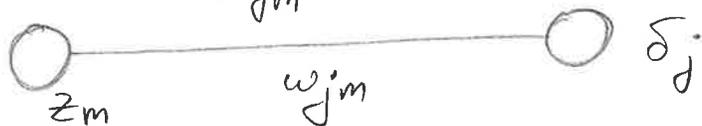$$\frac{\partial E}{\partial w_{jm}} = \delta_j \cdot z_m$$

$z_m$     $w_{jm}$     $\delta_j$    Fig. 3

— To compute derivatives, it suffices to:

(1) • compute values of $\delta_j$ for each node

(2) • also, save values $z_m$ during forward pass through Network

(will be multiplied w/ values of $\delta_j$ to get partials)

## How to compute the errors?

— For an output-layer unit (indexed by $\boxed{k}$),
  assume regression problem (i.e., output activation function is linear)
  assume LSE loss function:

(9)
$$\delta_k^{(\ell)} = \frac{\partial L}{\partial a_k^{(\ell)}} = \frac{\partial L}{\partial \hat{y}_k} = \frac{d\left(\frac{1}{2}(\hat{y}_k - y_k)^2\right)}{d\hat{y}_k} = \hat{y}_k - y_k$$

$k = 1 \div K$

— For a hidden-layer unit (indexed by $\boxed{j}$), use chain-rule again:

(10)
$$\delta_j^{(\ell)} = \frac{\partial L}{\partial a_j^{(\ell)}} = \sum_{m=0}^{M} \boxed{\frac{\partial L}{\partial a_m^{(\ell+1)}}} \cdot \boxed{\frac{\partial a_m^{(\ell+1)}}{\partial a_j^{(\ell)}}}$$

← see Figure. 2

$j = 0 \div M$

— all $\boxed{M}$ nodes that node $\boxed{j}$ sends connections to!

— Activation value of unit $\boxed{j}$ contributes only via its contribution to the activation value of those unit it is connected to! (i.e., $a_m^{(\ell+1)}$)

**1st term** usually non-linear (dependence between Loss and activation value of a unit on next layer)

**term #2** captures relationship between this activation $a_j^{(\ell)}$ and subsequent activation $a_m^{(\ell+1)}$ on next layer

Furthermore, we note:

term #1: $\dfrac{\partial L}{\partial a_m^{(\ell+1)}} \overset{\Delta}{=} \delta_m^{(\ell+1)}$ by definition $\qquad$ (11)

term #2: $\dfrac{\partial a_m^{(\ell+1)}}{\partial a_j^{(\ell)}} \overset{\text{by chain rule}}{=} \dfrac{\partial a_m^{(\ell+1)}}{\partial h(a_j^{(\ell)})} \cdot \dfrac{\partial h(a_j^{(\ell)})}{\partial a_j^{(\ell)}} = \omega_{mj}^{(\ell+1)} \cdot h'(a_j^{(\ell)})$ (12)

$$= \dfrac{d h(a_j^{(\ell)})}{d a_j^{(\ell)}} \overset{\Delta}{=} h'(a_j^{(\ell)})$$

$$= \dfrac{\partial a_m^{(\ell+1)}}{\partial z_j^{(\ell)}} = \omega_{mj}^{(\ell+1)}$$

Substitute then 2 terms in equation (**) from previous page to get:

(13) $\boxed{\delta_j^{(\ell)} = h'(a_j^{(\ell)}) \cdot \sum_{m=0}^{M} \omega_{mj}^{(\ell+1)} \cdot \delta_m^{(\ell+1)}}$ **: this is a key insight in Backpropagation !!!**

—— Value of "error" can be computed by "passing back" (i.e, eg backpropagating) the "errors" for nodes further up in the network !

- We **know values of errors** $\delta$ for the final/output layer
  — by forward calculations
- We can **recursively apply (13)** to backpropagate errors from output units toward input.

- We can then **use eq. (8) to estimate partial derivatives** of Loss function $L$ w.r.t. any weight !
  — those **derivatives can be used by opt. methods like GD to train all weights !!!**

*Loss Function — Network used for Logistic Regression* (6)

Log-Likelihood from Maximum Likelihood Estimation (MLE)

Assume probabilistic model, training pairs $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$, $y_i \in \{0, 1\}$

Let neural network w/ parameters $\theta$ produce a logit $z_i = z_\theta(x_i)$

and probability:

$$P_i \equiv P_\theta(y_i = 1 \mid x_i) = \sigma(z_i) = \frac{1}{1 + e^{-z_i}}$$

Under a Bernoulli model:

$$P(y_i \mid x_i; \theta) = \underbrace{P_i^{y_i}(1 - P_i)^{1 - y_i}}_{\text{compact version}} = \begin{cases} P_i, & \text{if } y_i = 1 \\ 1 - P_i, & \text{if } y_i = 0. \end{cases}$$

expanded PMF of
Bernoulli distribution.

Then, likelihood and log-likelihood (assuming of course iid):

$$\mathcal{L}(\theta) = \prod_{i=1}^{n} P_i^{y_i}(1 - P_i)^{1 - y_i}$$

$$\ell(\theta) = \log(\mathcal{L}(\theta)) = \sum_{i=1}^{n} \left[ y_i \cdot \log(P_i) + (1 - y_i) \log(1 - P_i) \right]$$

The average per-sample loss (divide by $\boxed{n}$) is the binary cross-entropy:

$$J(\theta) \equiv \boxed{BCE(y, p) = -\frac{1}{n} \sum_{i=1}^{n} \left[ y_i \cdot \log P_i + (1 - y_i) \cdot \log(1 - P_i) \right]}$$

↑
popular
notation

So, BCE is exactly the NLL of a Bernoulli model with

$$p = \sigma(z) \, !$$

Minimizing BCE is equivalent to maximum likelihood for binary classification!