# Multi-Layer Perceptron
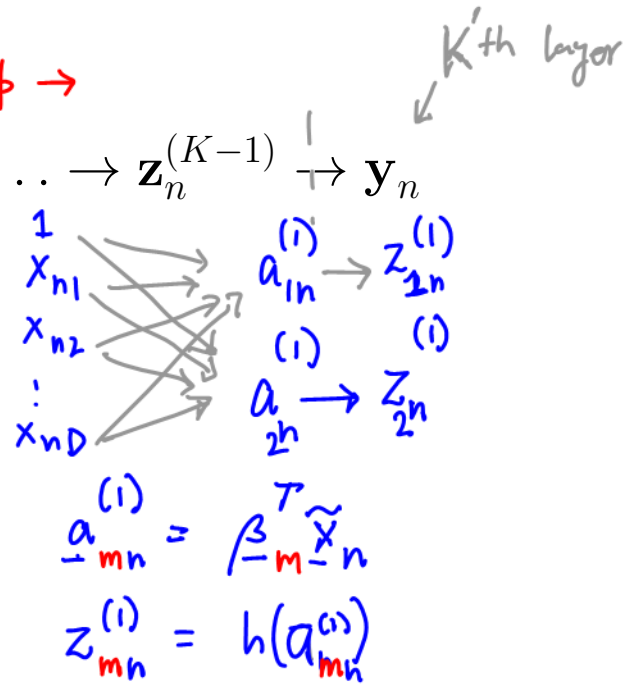
Mohammad Emtiyaz Khan
EPFL

Dec. 3, 2015

# Multi-Layer Perceptron (MLP)

This is also known as feed-forward neural network and can be represented graphically as follows:

$$\mathbf{x}_n \to \mathbf{a}_n^{(1)} \to \mathbf{z}_n^{(1)} \to \mathbf{a}_n^{(2)} \to \mathbf{z}_n^{(2)} \to \ldots \to \mathbf{z}_n^{(K-1)} \to \mathbf{y}_n$$
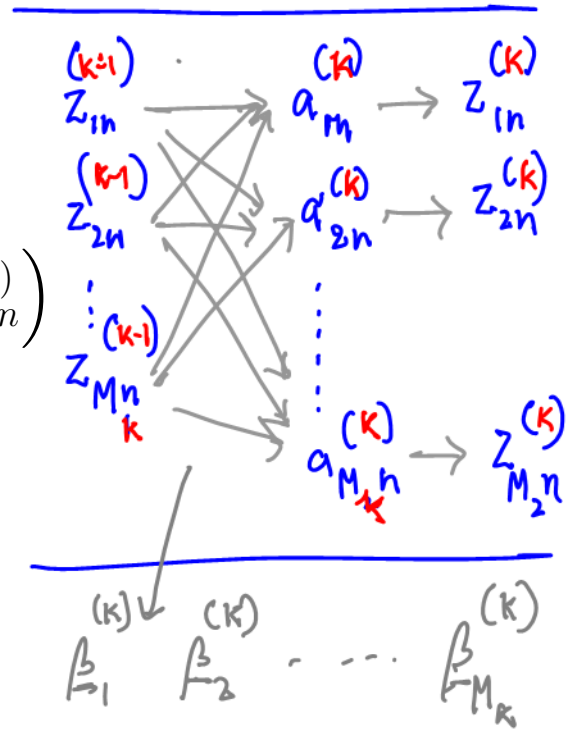
where $\{y_n, \mathbf{x}_n\}$ is the $n$'th input-output pair, $\mathbf{z}_n^{(k)}$ is the $k$'th hidden vector, $\mathbf{a}_n^{(k)}$ is the corresponding activation. There are a total of $K$ layers.
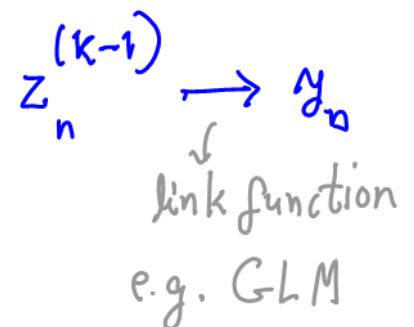
For the $k$'th layer, we obtain the $m$'th activation $a_{mn}^{(k)}$ and the corresponding hidden variable $z_{mn}^{(k)}$, as shown below:

$$a_{mn}^{(k)} = \left(\boldsymbol{\beta}_m^{(k)}\right)^T \mathbf{z}_n^{(k-1)}, \quad z_{mn}^{(k)} = h\left(a_{mn}^{(k)}\right)$$

where $\mathbf{z}_n^{(k-1)}$ is the hidden vector for the previous layer. For the first layer, we set $\mathbf{z}_n^{(0)} = \mathbf{x}_n$. For the last layer, we use a link function to map $\mathbf{z}_n^{(K-1)}$ to the output $\mathbf{y}_n$.

Note that a 1-Layer MLP is simply a generalization of linear/logistic regression.

$$(y_n, \underline{x}_n)$$

$$\text{deep} \to$$

$$K\text{'th layer}$$

$$a_{mn}^{(1)} = \boldsymbol{\beta}_m^T \tilde{\underline{x}}_n$$

$$z_{mn}^{(1)} = h\left(a_{mn}^{(1)}\right)$$

$$z_n^{(k-1)} \longrightarrow y_n$$

link function

e.g. GLM

1

Defining $\mathbf{B}^{(k)}$ as a matrix with rows $(\boldsymbol{\beta}_m^{(k)})^T$, we can express the computation of activation and hidden vectors as follows:
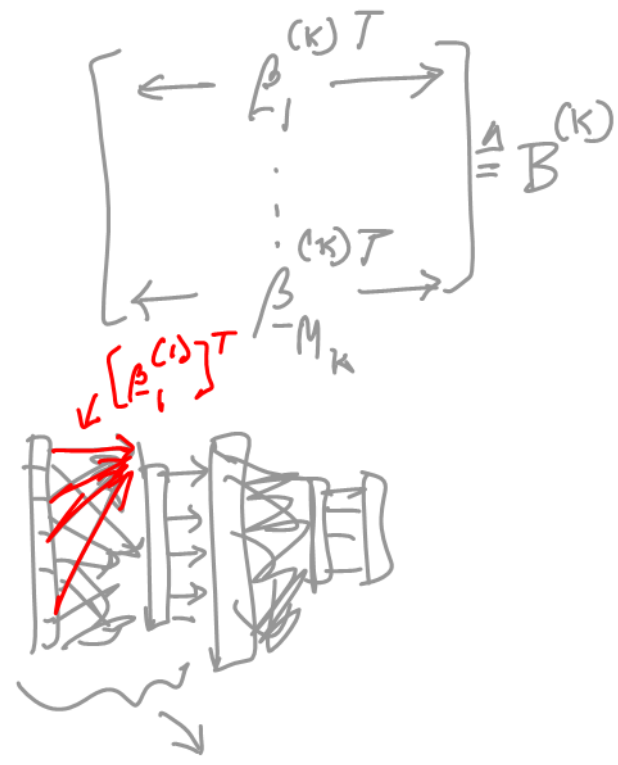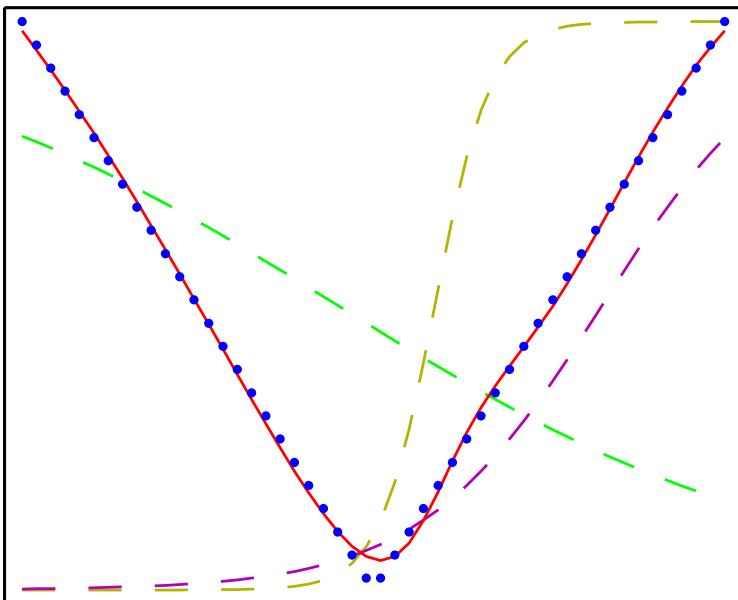
$$\mathbf{a}_n^{(k)} = \mathbf{B}^{(k)}\mathbf{z}_n^{(k-1)}, \quad \mathbf{z}_n^{(k)} = h\left(\mathbf{a}_n^{(k)}\right)$$

In a more compact notation, we can express the input-output relationship as follows:

$$\hat{y}_n = g((\boldsymbol{\beta}^{(K-1)})^T * h(\mathbf{B}^{(K-2)} * h(* \ldots * h(\mathbf{B}^{(1)} * \mathbf{x}_n)))),$$

where $g$ is an appropriate link function to match the output.

An illustration below shows reconstruction of the function $|x|$ at $N = 50$ data points sampled at the blue dots. The trained network has 2 layers and 3-hidden variables with $tanh()$ activation function.

$$\left[ \begin{matrix} \longleftarrow & \beta_1^{(k)\,T} & \longrightarrow \\ & \vdots & \\ \longleftarrow & \beta_{-M_k}^{(k)\,T} & \longrightarrow \end{matrix} \right] \triangleq \mathbf{B}^{(k)}$$

$$\swarrow \left[\beta_1^{(1)}\right]^T$$

$$\sum_{n=1}^{N} \left( y_n - \hat{y}_n \right)^2$$

$$f\left(B^{(1)}, \ldots, {}^{(k-1)}\right)$$

$$\frac{\partial f}{\partial B^{(k)}}$$



2

# Optimization and Back-propagation

We can learn parameters **B** using stochastic gradient-descent.

Gradient computation can be complicated due to the *deep* structure of the network. We can use back-propagation to simplify the computation. The key-idea is to express the derivatives in terms of activations $\mathbf{a}_n^{(k)}$ and hidden variables $\mathbf{z}_n^{(k)}$ using the chain rule. Below is the outline of the algorithm:

$$\ell = \sum_{n=1}^{N} \left( y_n - f(\underline{x}_n) \right)^2$$

$$\sum_{n=1}^{N} \left( y_n - g\left( B^k * h\left( B^{k-1} * \cdots (\underline{x}_n) \right) \right) \right)^2$$

$$K \times H \times 256 \times 256$$

$$1,00,0,000 \qquad 10,000 \leftarrow$$

$$B = \{ B^{(1)}, B^{(2)}, \cdots B^{(k)} \}$$

$$\frac{\partial \ell}{\partial B} = \sum_{n=1}^{N} \begin{bmatrix} 10,600 \end{bmatrix}$$

$$\text{Use the chain-rule} \longrightarrow \widehat{\frac{\partial \ell}{\partial B}} = \begin{bmatrix} 10,000 \end{bmatrix}$$

$$B_{i+1} = B_i - \delta_k \left. \frac{\widehat{\partial \ell}}{\partial B} \right|_{B = B_i}$$

$$\left. \frac{\partial \ell}{\partial B^{(k-1)}} \right|_{B^{(k-1)} = B_i^{(k-1)}}$$

$$\downarrow \frac{\partial \ell}{\partial B^{(k-2)}}$$

$$f(\underline{x}_{n_i})$$

$$B_i \uparrow$$

$$X_{n_i}$$

$$\frac{\partial \ell}{\partial B^{(1)}}$$



**Forward Pass**

$a^{(2)}$

$z_1^{(1)}$ $z_2^{(1)}$ $z_3^{(1)}$

$a_1^{(1)}$ $a_2^{(1)}$ $a_3^{(1)}$

$x_1$ $x_2$

**Backward Pass**

$r^{(2)}$

$r^{(2)} w_1^{(2)}$ $r^{(2)} w_2^{(2)}$ $r^{(2)} w_3^{(2)}$

$* g'\left(a_1^{(1)}\right)$ $* g'\left(a_2^{(1)}\right)$ $* g'\left(a_3^{(1)}\right)$
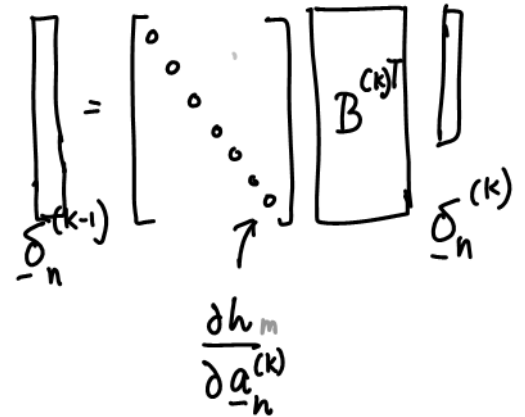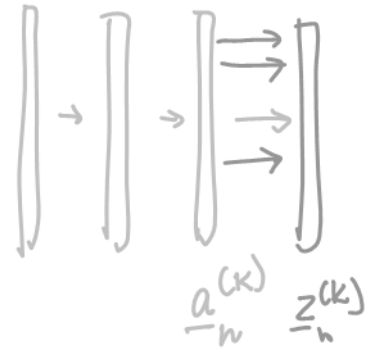
$r_1^{(1)}$ $r_2^{(1)}$ $r_3^{(1)}$

Step 1: Compute $\mathbf{a}_n^{(k)}$ and $\mathbf{z}_n^{(k)}$ using forward propagation.

Step 2: Compute $\boldsymbol{\delta}_n^{(k)} := \partial\mathcal{L}/\partial\mathbf{a}_n^{(k)}$ using backward propagation:

$$\boldsymbol{\delta}_n^{(k-1)} = \text{diag}\left[\mathbf{h}'(\mathbf{a}_n^{(k)})\right]\left(\mathbf{B}^{(k)}\right)^T\boldsymbol{\delta}_n^{(k)}$$

Step 3: Compute $\partial\mathcal{L}/\partial\mathbf{B}^{(k)}$ using the above derivatives.

$$\frac{\partial\mathcal{L}}{\partial\mathbf{B}^{(k)}} = \sum_n \boldsymbol{\delta}_n^{(k)}\left(\mathbf{z}_n^{(k)}\right)^T$$

## Tricks

Obtaining a good generalization error with neural networks and avoiding overfitting requires a lot of hacks and tricks. A good summary of these are given in Bottou's paper "Stochastic gradient tricks". In addition, initialization seems to play a huge role in improving the performance. See the following paper "On the importance of initialization and momentum in deep learning" by Ilya Sutskever et. al.