

Multi-Layer Perceptron

Mohammad Emtiyaz Khan
EPFL

Dec. 3, 2015



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

Multi-Layer Perceptron (MLP)

This is also known as [feed-forward neural network](#) and can be represented graphically as follows:

$$\mathbf{x}_n \rightarrow \mathbf{a}_n^{(1)} \rightarrow \mathbf{z}_n^{(1)} \rightarrow \mathbf{a}_n^{(2)} \rightarrow \mathbf{z}_n^{(2)} \rightarrow \dots \rightarrow \mathbf{z}_n^{(K-1)} \rightarrow \mathbf{y}_n$$

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.

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:

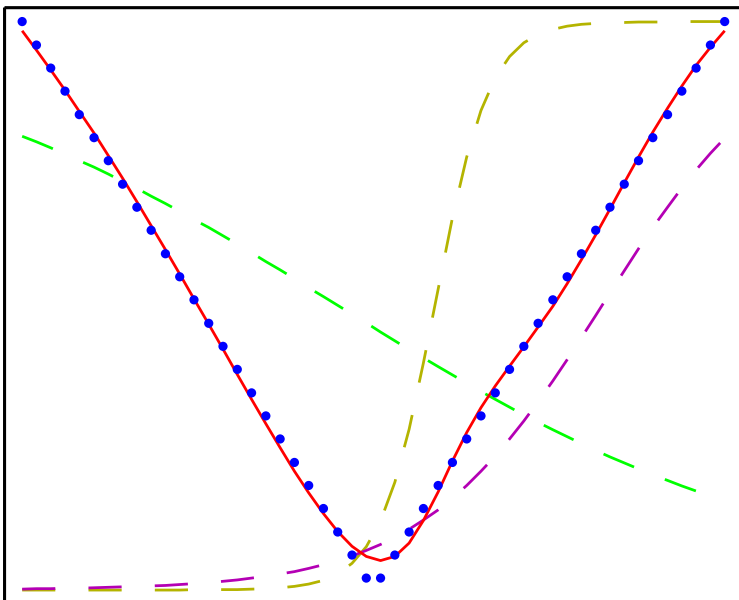
$$\mathbf{a}_n^{(k)} = \mathbf{B}^{(k)} \mathbf{z}_n^{(k-1)}, \quad \mathbf{z}_n^{(k)} = h\left(\mathbf{a}_n^{(k)}\right)$$

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

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

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

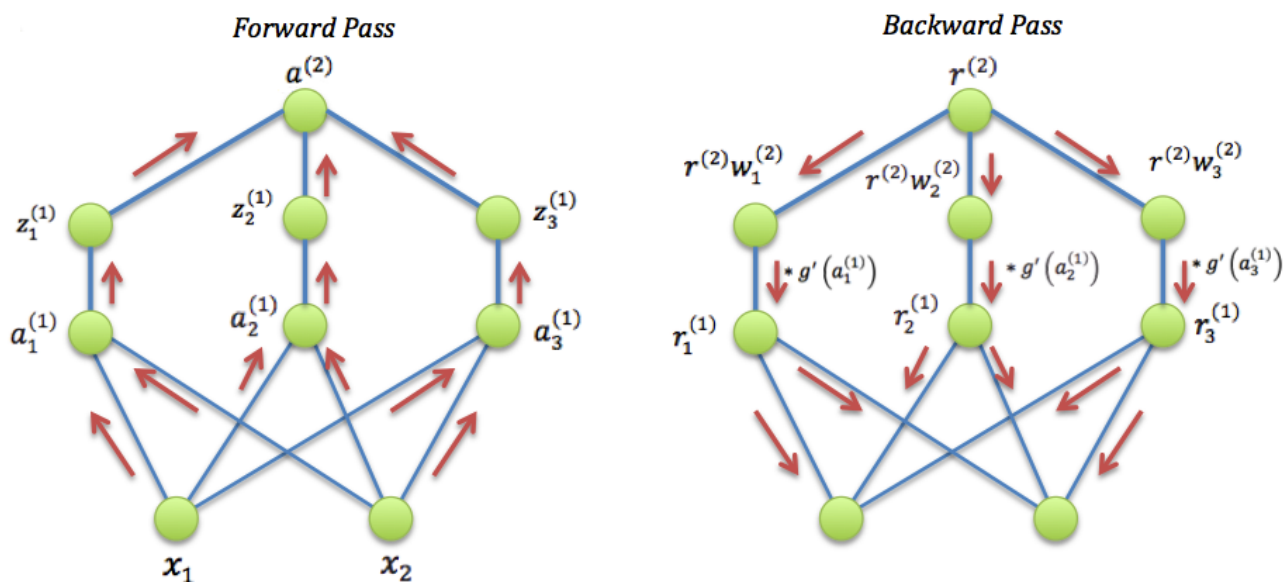
An illustration below shows reconstruction of the function $|x|$ at $N = 50$ data points sampled at the blue dots. The trained network has 2 layers and 3-hidden variables with $\tanh()$ activation function.



Optimization and Back-propagation

We can learn parameters \mathbf{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:



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.