where $w_k$ are the singular values computed by SVD.

## 3.4  Sigma-clipping

As one last set of problems explored here, consider the need to fit the data *asymmetrically*; for example, we may want to fit a baseline, or an envelope of a signal.

Fortunately, there is no new formalism that needs to be developed; we will use linear least squares explained above in a special way.

1. Given data points $(x_i, y_i)$, fit a model to all data.

2. Compute the residuals $\Delta_i = y_i - y(x_i|\boldsymbol{c})$ and the variance of the residuals, $\sigma_0$.

3. Given upper and lower clipping limits, $\eta_{\mathrm{hi}}$ and $\eta_{\mathrm{lo}}$, remove all points that are above $\eta_{\mathrm{hi}}\sigma_0$ and below $\eta_{\mathrm{lo}}\sigma_0$.

4. Repeat the fit on the reduced number of data.

5. Continue until no further points are removed from the data.

This way, by controlling $\eta_{\mathrm{lo}}$ and $\eta_{\mathrm{hi}}$, we control the part of the data-set that the model will fit.

# 4  Artificial neural networks

A range of problems in computational physics can be very difficult or even impossible to describe analytically, or the forward models take such a long time that any attempts at modeling are rendered impractical. In those cases we can consider an alternative: artificial neural networks.

An artificial neural network is a simple construct: it is a stack of interconnected *layers* (cf. Fig. 1). Each layer is an array of processing elements called *units*. These units propagate the signal between layers by *weighted connections*. The units perform non-linear *mapping* of input data to output parameters.

The input to the network is a measured or observed dataset – an array of sampled data points. These data are the stimulus of the *input* layer: each unit on the layer acquires a value of the given input array element. Once the input layer is populated, propagation to the *hidden* layer begins. The stimulus of the given unit on the hidden layer is a weighted sum of outputs from units on the current layer. This stimulus is then passed through a non-linear *activation function* that determines the extent of stimulation of the given unit. There is a connection from every unit on the input layer to every unit on the hidden layer, and each of these connections has a corresponding weight assigned to it. Once the signal has been propagated to the hidden layer, the propagation continues to the *output* layer: each unit acquires a value that is a weighted sum of outputs from units on the hidden layer, passed through the activation function.

Figure 1: ADD FIGURE! The topology of a three-layer Artificial Neural Network (ANN). Processing units (nodes) on the input layer are populated with input observations. A weighted sum of values on the input layer, $y_j = \sum_k w_{jk}^H i_k$, stimulates each unit, $j = 1 \ldots l$, on the hidden layer. The amount of stimulus is determined by the non-linear activation function $A_f^H$, such as a sigmoid function, $A_f^H(y_j) = 1/[1 + \exp(-(y_j - \mu)/\tau)]$, with parameters $\mu$ and $\tau$ chosen in such a way that $A_f^H$ is mapped onto the $[-1, 1]$ interval. Processing units on the hidden layer are then populated with $h_j = A_f^H(y_j)$. An analogous propagation takes place between the hidden layer and the output layer, using the same type of the activation function. The output layer of a trained network then contains the network's best guess at descriptive parameters of the data that have populated the input layer. The ANN is thus a *non-linear mapping* from observations to the physical parameters of the observed system.

The network thus *maps* its input to its output by propagating the stimulus via weighted connections that are passed through activation functions. Given the layer-to-layer connection weights, the propagation is basically a matter of summation and multiplication. The power comes from the ability of networks to provide *non-linear* mapping between their input and their output by using non-linear activation functions. The most commonly used activation functions are sigmoid curves, i.e. functions of the type $f(x) = 1/[1 + \exp(-(x - \mu)/\tau)]$.

The goal of ANNs is to have output layer values that have some representative significance of the input array. Since the output values depend on the connection weights, the task is to determine these weights. This is where *back-propagation* comes into play. Assume we have a sample of several thousand *exemplars* (input arrays for which we know the corresponding parameters) – obtained, for example, by computing theoretical models, or by using real data with reliable model solutions. For each exemplar we perform forward-propagation and compare the results output by the network with the true ones. We then *modify* the weights so that the discrepancy between the results is reduced for the whole sample. This iterative procedure is called the *training* or *machine learning* phase. It is the only computation-intensive block and it is performed only once. Once the weights are determined and the network reliably reproduces the expected results, the network is ready to *recognize* input never seen before. Hundreds of thousands of inputs can subsequently be processed in a matter of seconds on a single CPU.

The network described here, a basic three-layer back-propagating network (BPN), is remarkably robust for a diversity of non-linear problems. Different network topologies, such as multiple hidden layers, as well as more complicated connection strategies, advanced training approaches and other variations in general do not bring significant improvements to the basic model. For a thorough discussion on neural networks please refer to authoritative books such as Freeman & Skapura (1991).

## 4.1   Back-propagation

Back-propagation, as stated above, is the modification of weights to minimize the residuals between exemplar outputs and network-produced outputs on the output layer. As the mapping from input to output layer is multi-dimensional and non-linear, we resort to steepest descent.

The net input to the $j$-th unit on the hidden layer is:

$$\xi_{pj}^h = \sum_k w_{jk}^h i_k.$$

That net input is "activated" by the activation function $A_h(\xi_{pj}^h)$. Similarly for the output layer:

$$\xi_{pi}^o = \sum_j w_{ij} h_j,$$

activated by $A_o(\xi_{pi}^o)$. We want to minimize:

$$E_p = \frac{1}{2} \sum_p \delta_{pi}^2, \quad \text{where } \delta_{pi} = y_{pi} - o_{pi}.$$

To do that, we compute the (negative) gradient of $E_p$:

$$-\frac{\partial E_p}{\partial w_{ij}} = -\delta_{pi} \frac{\partial A_o(\xi_{pi}^o)}{\partial \xi_{pi}^o} \frac{\partial \xi_{pi}^o}{\partial w_{ij}^o}. \tag{45}$$

Given the functional form of $A_o(\xi_{pi}^o)$, the first derivative is simple to calculate; two common examples are linear and sigmoid:

$$A_o(\xi_{pi}^o) = \xi_{pi}^o, \quad \frac{\partial A_o(\xi_{pi}^o)}{\partial \xi_{pi}^o} = 1;$$

$$A_o(\xi_{pi}^o) = \left(1 + e^{-\xi_{pi}^o}\right)^{-1}, \quad \frac{\partial A_o(\xi_{pi}^o)}{\partial \xi_{pi}^o} = \xi_{pi}^o(1 - \xi_{pi}^o) \equiv o_{pi}(1 - o_{pi}).$$

The second derivative is also simple to calculate:

$$\frac{\partial \xi_{pi}^o}{\partial w_{ij}^o} = \frac{\partial}{\partial w_{ij}^o} \sum_j w_{ij} h_j = h_j.$$

Thus,

$$-\frac{\partial E_p}{\partial w_{ij}} = -\delta_{pi} A_o'(\xi_{pi}^o) h_j \equiv \delta_{pi}^o h_j.$$

Here we introduced $\delta_{pi}^o \equiv \delta_{pi} A_o'(\xi_{pi}^o)$. Then we correct the weights by taking a step along the gradient:

$$w_{ij}^o \mapsto w_{ij}^o + \eta \delta_{pi}^o h_j. \tag{46}$$

We do the exact same procedure for adjusting weights from the input layer to the hidden layer.

To train the neural network, we thus:

1. apply the input vector $\boldsymbol{i} = (i_1, \ldots, i_N)$ to the $N$ input units;

2. calculate the net input values to the hidden layer units as:

$$h_j = A_h \left( \sum_k w_{jk}^h i_k \right), \quad \text{where } A_h(x) = \left( 1 - e^{-(x-\mu)/\tau} \right)^{-1}; \quad (47)$$

3. calculate the net input values to the output layer as:

$$o_i = A_o \left( \sum_j w_{ij}^o h_j \right), \quad \text{where } A_o(x) = \left( 1 - e^{-(x-\mu)/\tau} \right)^{-1}; \quad (48)$$

4. calculate the error terms for the output units:

$$\delta_i^o = (y_i - o_i) A_o' \left( \sum_j w_{ij}^o h_j \right), \quad (49)$$

where $y_i$ are known values of parameters;

5. calculate the error terms for the hidden units:

$$\delta_j^h = A_h' \left( \sum_k w_{jk}^h i_k \right) \sum_i \delta_i^o w_{ij}^o; \quad (50)$$

6. update weights on the output layer:

$$w_{ij}^o \mapsto w_{ij}^o + \eta \delta_i^o A_o \left( \sum_j w_{ij}^o h_j \right), \quad (51)$$

where $\eta$ is the learning rate parameter;

7. update weights on the hidden layer:

$$w_{jk}^h \mapsto w_{jk}^h + \eta \delta_j^h i_j; \quad (52)$$

8. calculate the overall error term:

$$E_p = \frac{1}{2} \sum_i \delta_i^2 \quad \text{for the } p\text{-th exemplar}; \quad (53)$$

9. repeat all steps for all $P$ exemplars.

The idea is to keep training the network for as long as $E_p$ are decreasing and stop once they dip below a certain threshold.

So far we only talked about neural networks used for regression; another frequent application is classification.