

## I. THEORY

### A. Neural network

A neural network consists of neurons, also called nodes, in a network which are inspired by the neurons in the brain (?). The node takes input from the neighbouring nodes and uses an activation function  $f$  to calculate the input nodes with the weights  $w$  that describes the connection between the nodes. The output  $y$  that is connected to a set of nodes  $x_i$  with weights  $w_i$  is gives an

$$y = f\left(\sum_{i=1}^n w_i x_i\right) = f(u) \quad (1)$$

and is illustrated in Figure 1a. A neural network consist of layers with connected neurons, or also called nodes.

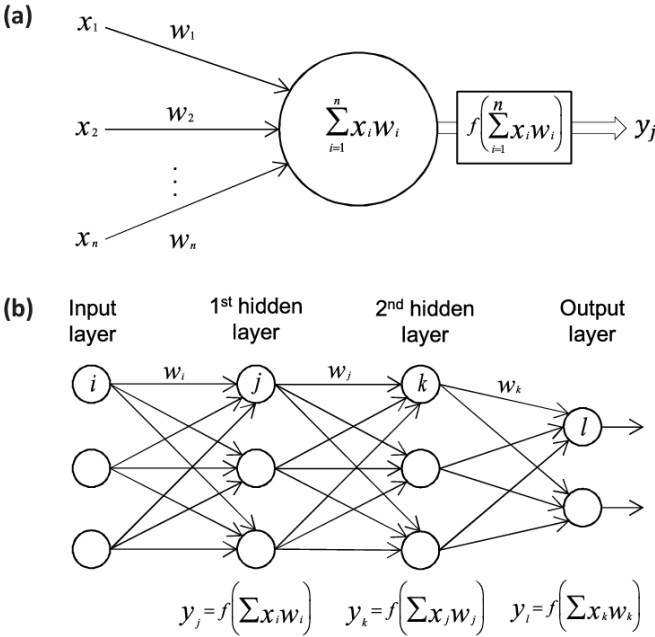


FIG. 1: Neural network, from (?)

The output  $y$  of a node is given via the activation function  $f$  as

$$y = f\left(\sum_{i=1}^n w_i x_i + b_i\right) = f(z), \quad (2)$$

where  $x_i$  is the input from the neurons in the preciding layers, and  $b_i$  is called the *bias* which is included in the case of a zero weight  $w_i$  or input  $x_i$  in the activation function. In e.g the first layer of a neural network we calculate the weighted sum  $z_i^1$  for each node in the layer over the input coordinates  $x_j$  as

$$z_i^1 = \sum_{j=1}^M w_{ij}^1 x_j + b_i^1, \quad (3)$$

where  $M$  denotes all the possible inputs to a given node  $i$  in the first layer. We call the layer *fully connected* if all the nodes in the layer takes all the nodes from the previous layer as input, and  $M$  is then equal to the number of nodes in the previous layer. The activation function  $f_i$  for each node  $i$  takes the value of the activation function in Eq. (3) and gives us the output  $y_i^1$  of all neurons in layer 1 as

$$y_i^1 = f(z_i^1) = f\left(\sum_{j=1}^M w_{ij}^1 x_j + b_i^1\right), \quad (4)$$

with the assumption that all the nodes in the same layer has the same activation function, and we therefore drop the subscript  $i$  on  $f_i$ . If the layers has different activation functions we can denote it with a superscript  $l$  for a layer  $l$ . The output of the node  $i$  for layer  $l$  therefore becomes

$$y_i^l = f(z_i^l) = f\left(\sum_{j=1}^{N_{l-1}} w_{ij}^l y_j^{l-1} + b_i^l\right), \quad (5)$$

where  $N_{l-1}$  is the number of nodes in the previous layer if the layer is fully connected. With Eq. (5) we can calculate the output  $y_i^l$  from layer  $l$  from the output  $y_i^{l-1}$  from the previous layer  $l-1$ . Given the input  $x_n$  from the input layer with  $n$  nodes, the output from layer  $l+1$  can therefore be calculated as

$$y_i^{l+1} = f^{l+1}\left[\sum_{j=i}^{N_l} w_{ij}^{l+1} f^l\left(\sum_{k=i}^{N_{l-1}} w_{jk}^l \left(\dots f^1\left(\sum_{m=i}^{N_0} w_{mn}^1 x_n + b_m^1\right) \dots\right) + b_j^l\right) + b_i^{l+1}\right]. \quad (6)$$

Eq. (6) shows that the MPL is nothing more than an analytic function which is given by the real valued map  $\hat{x} \in \mathbb{R}^n \rightarrow \hat{y} \in \mathbb{R}^m$ . This concludes the feed forward method. Now we go over to the famous back propagation algorithm (cite) which is used to update the weights and biases in the neural network. To summarize the weights  $w_{ij}^l$  and biases  $b_j$  are updated according to

$$w_{jk}^l \leftarrow w_{jk}^l - \eta \delta_j^l y_k^{l-1} \quad (7)$$

$$b_j^l \leftarrow b_j^l - \eta \delta_j^l, \quad (8)$$

where  $\delta_j^l$  is the error that is calculated from the derivative of the cost function  $\mathcal{C}$  and  $\eta$  is the learning rate. A more detailed description of the back propagation algorithm can be found in (?).

### B. Discovery statistics and the profile likelihood ratio in particle physics

In particle physics we often test a so called *background hypythis* that test the physics we know against the *signal + background hypotisis* that test unknown physics

that could possibly be discovered by rejecting the background hypothesis. One way of doing this is to calculate the discovery significance  $Z$  which counts the number of standard deviations and is the square of some test statistics  $q_0$ . In a search region where there are  $b$  background events and  $s$  signal events the discovery significance is given as

$$Z = \sqrt{2(s+b) \ln \left(1 + \frac{s}{b}\right)} - s, \quad (9)$$

and is often used in particle physics. Another way of using the test statistics is to use the *profile likelihood ratio* and the signal strength  $\mu$  which measures how strong the signal is. To test a hypothesis value for the signal strength  $\mu$  we use the *profile likelihood ratio* (?)

$$\lambda(\mu) = \frac{L(\mu, \hat{\theta})}{L(\hat{\mu}, \hat{\theta})}, \quad (10)$$

where in the numerator the value of  $\theta$  that maximizes  $L$  for a specific  $\mu$  is denoted by the quantity  $\hat{\theta}$ , and is therefore the conditional maximum-likelihood (ML) estimator that estimates  $\theta$ , and therefore a function of  $\mu$ . The denominator is the unconditional likelihood function where  $\hat{\mu}$  and  $\hat{\theta}$  are their ML estimators.

By Wilks' theorem we can use the test statistics

$$q_0 = \begin{cases} -2 \ln \lambda(0), & \text{if } n > b, \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

The appropriate likelihood ratio in searches for new particles is the ratio of the likelihood for the signal+background hypothesis and the likelihood for the background hypothesis and is given generally as

$$Q = \frac{L_{s+b}(\mu)}{L_b} \quad (12)$$

This likelihood ratio can be expressed in terms of the background only likelihood function  $L_b$  and signal only density function  $L_s$  and is given as (?)

$$Q = e^{-s_{exp}} \prod_{i=1}^n \left(1 + \frac{s_{exp} L_s(x_i)}{b_{exp} L_b(x_i)}\right), \quad (13)$$

where  $b_{exp}$  and  $s_{exp}$  are the expected background and signal events, and the product is taken over a given experimental result  $x_i$  for each event  $i$  over the total number of events  $n$ . We insert the signal strength parameter  $\mu$  in front of the expected number of signal events  $s_{exp}$ , which gives us

$$Q(\mu) = e^{-\mu s_{exp}} \prod_{i=1}^n \left(1 + \frac{\mu s_{exp} L_s(x_i)}{b_{exp} L_b(x_i)}\right). \quad (14)$$

The profile log-likelihood ratio  $t(\mu)$  is therefore equal to

$$t(\mu) = -2 \ln Q(\mu) = \mu s_{exp} - \sum_{i=1}^n \left(1 + \frac{\mu s_{exp} L_s(x_i)}{b_{exp} L_b(x_i)}\right), \quad (15)$$

when expressed in terms of  $L_s$  and  $L_b$ . In the profile likelihood ratio in Eq. (10) we see that the estimator value that maximizes the signal+background likelihood function is  $\hat{\mu} = 1$  in the unconditional likelihood function in the denominator in Eq. (10). Using this the derivation in Eq. .. in Appendix gives us the relation

$$\begin{aligned} -2 \ln \lambda(\mu) &= -2 \ln Q(\mu) + 2 \ln Q(1) \\ &= t(\mu) - t(1), \end{aligned} \quad (16)$$

If the estimated value  $\hat{\mu}$  that maximises an *approximated* signal+background likelihood function is not equal to the expected estimate  $\hat{\mu} = 1$ , as used to derive Eq. (16), we use the estimated value  $\hat{\mu}$  instead in the relation above. This gives us the same relation as above, but where we replace 1 by  $\hat{\mu}$  in the last term

$$\begin{aligned} -2 \ln \lambda(\mu) &= -2 \ln Q(\mu) + 2 \ln Q(\hat{\mu}) \\ &= t(\mu) - t(\hat{\mu}), \end{aligned} \quad (17)$$

which is the more general form. This gives us the test statistic

$$q_0 = \begin{cases} t(0) - t(\hat{\mu}), & \text{if } n > b_{exp}, \\ 0, & \text{otherwise,} \end{cases} \quad (18)$$

and the discovery significance  $Z$  is then given as

$$Z = \sqrt{q_0} = \begin{cases} \sqrt{t(0) - t(\hat{\mu})}, & \text{if } n > b_{exp}, \\ 0, & \text{otherwise.} \end{cases} \quad (19)$$

### C. Likelihood estimation

The likelihood ratio function can be estimated by a neural network with a method developed by G. V. Moustakides and K. Basioti (?) which transforms the output of the neural network to a likelihood ratio by solving an optimization problem. For a random vector  $X$  with the probability function  $f_0(X)$  we want to optimize the following loss function

$$\mathcal{J}(u) = E_0[\phi(u(X)) + r(X)\psi(u(X))], \quad (20)$$

where  $E_0[\cdot]$  is the expectation value w.r.t.  $f_0(X)$ , and  $\phi(z)$  and  $\psi(z)$  are scalar functions of the scalar  $z$ , and  $u(X)$  and  $r(X)$  are scalar functions of  $X$  where  $r(X)$  takes values in a *known* interval  $\mathcal{I}_r$ . The optimization problem we want to solve is

$$\min_{u(X)} \mathcal{J}(u) = \min_{u(X)} E_0[\phi(u(X)) + r(X)\psi(u(X))], \quad (21)$$

such that when we have a known scalar function  $\omega(r)$  of the scalar  $r$ , we want to determine the two functions  $\phi(z)$  and  $\theta(z)$  such that the global minimizer in Eq. (??) is equal to

$$u(X) = \omega(r(X)). \quad (22)$$

In (?) they show that the necessary condition to achieve the optimization that satisfies this condition is that

$$\phi'(\omega(r)) + r\psi'(\omega(r)) = 0. \quad (23)$$

If  $\omega(r)$  is strictly monotone we also have the result

$$\phi'(\omega(r)) + \omega^{-1}(z)\psi'(\omega(r)) = 0, \quad (24)$$

for  $z \in \omega(\mathcal{I}_r)$  where  $\omega^{-1}(z)$  is the inverse of  $r$ . They further show that when we chose a function  $\rho(z)$  with the condition  $\rho(z) < 0$  for all  $z \in \omega(\mathcal{I}_r)$  and define  $\phi(z)$  and  $\theta(z)$  as

$$\psi'(z) = \rho(z), \text{ and } \phi'(z) = \omega^{-1}(z)\rho(z), \quad (25)$$

then the minima of  $\phi(z) + r\theta(z)$  in Eq. (21) for  $z \in \omega(\mathcal{I}_r)$  is uniquely located at  $z = \omega(r)$ . This is what we wanted to achieve in Eq. (22). When we chose

$$\omega(z) = \frac{z}{z+1} \text{ and } \rho(z) = -\frac{1}{z}, \quad (26)$$

and use Eq. (25) we get the following expressions for  $\phi(z)$  and  $\theta(z)$ ,

$$\phi(z) = -\log(1-z) \text{ and } \psi(z) = -\log(z). \quad (27)$$

This combination of  $\phi(z)$  and  $\theta(z)$  gives us the binary cross-entropy method and is among the most popular methods for classification problems. When we define the scalar function  $r(X)$  in Eq. (20) to be the likelihood ratio

$$r(X) = \frac{f_1(X)}{f_0(X)}, \quad (28)$$

of two probability distributions  $f_0(X)$  and  $f_1(X)$  the loss function  $\mathcal{J}(u)$  in Eq. (20) becomes

$$\begin{aligned} \mathcal{J}(u) &= E_0[\phi(u(X)) + r(X)\psi(u(X))] \\ &= E_0[\phi(u(X))] + E_1[\psi(u(X))], \end{aligned} \quad (29)$$

where the expectation values  $E_0$  and  $E_1$  are w.r.t the probability functions  $f_0(X)$  and  $f_1(X)$  respectively. We are now interested in finding a way of using a neural network to estimate the transformation  $\omega(\frac{f_1(X)}{f_0(X)})$  when we

have two data samples  $\{X_1^0, \dots, X_{n_0}^0\}$  and  $\{X_1^1, \dots, X_{n_1}^1\}$  that are sampled from the *unknown* probability densities  $f_0(X)$  and  $f_1(X)$  respectively, and when  $\omega(r)$  is a *known* scalar function. Two approximations are needed to do this. For the first one we replace the statistical expectation values in Eq. (29) with the averages over the available data, while for the second we replace the function  $u(X)$  with the output of a neural network  $u(X, \theta)$ , where  $\theta$  summarizes the network parameters. The first approximation requires a dataset with a large sample size, while the second approximation requires that the neural network is sufficiently rich such that it can approximate any non-linear function. When applying the two approximations above, the loss function in Eq. (29) is approximated to

$$\mathcal{J}(u) \approx \hat{\mathcal{J}}(\theta) = \frac{1}{n_0} \sum_{i=1}^{n_0} \phi(u(X_i^0, \theta)) + \frac{1}{n_1} \sum_{i=1}^{n_1} \phi(u(X_i^1, \theta)) \quad (30)$$

The optimization problem is now transformed into an optimization problem over the neural network parameters  $\theta$  as

$$\begin{aligned} \min_{u(X)} \mathcal{J}(u) &\approx \min_{\theta} \hat{\mathcal{J}}(\theta) \\ &= \min_{\theta} \left\{ \frac{1}{n_0} \sum_{i=1}^{n_0} \phi(u(X_i^0, \theta)) + \frac{1}{n_1} \sum_{i=1}^{n_1} \phi(u(X_i^1, \theta)) \right\}, \end{aligned}$$

and when optimized the optimal parameters  $\theta_o$  will give a neural network that has approximated the function  $\omega(r)$ , that is

$$u(X, \theta_o) \approx \omega(r(X)), \quad (31)$$

provided that it does not reach a local minima.

#### D. Estimating the profile log-likelihood ratio with neural network

We can use the method of estimating a likelihood ratio described in the previous section when discovering new particles by estimating the profile log-likelihood ratio in Eq. (10). When we are interested to test the signal+background hypothesis against the background only hypothesis, which is given in the likelihood ratio function in Eq. (15), the likelihood ratio  $L_s/L_b$  is optimal to estimate when we have two datasets  $X_s$  and  $X_b$  which contains the signal and background data respectively