

Mass reconstruction with a neural network

L. Lönnblad¹

Deutsches Elektronen-Synchrotron DESY, Notkestraße 85, W-2000 Hamburg 52, FRG

C. Peterson² and T Rönvaldsson³

Department of Theoretical Physics, University of Lund, Sölvegatan 14A, S-22362 Lund, Sweden

Received 20 November 1991; revised manuscript received 6 January 1992

A feed-forward neural network method is developed for reconstructing the invariant mass of hadronic jets appearing in a calorimeter. The approach is illustrated in $W \rightarrow q\bar{q}$, where W-bosons are produced in $p\bar{p}$ reactions at SPS collider energies. The neural network method yields results that are superior to conventional methods. This neural network application differs from the classification ones in the sense that an analog number (the mass) is computed by the network, rather than a binary decision being made. As a by-product our application clearly demonstrates the need for using “intelligent” variables in instances when the amount of training instances is limited.

1. Introduction

Artificial neural network (ANN) have shown great promise for pattern recognition problems in particle physics. In particular, feed-forward classifiers have been used in off-line jet data analysis for quark/gluon separation and b-quark identification with remarkable success [1–5]. The principle is straightforward; a set of feature (classifier) units are parameterized as functions of kinematical data in a non-linear way by means of so-called sigmoidal functions (tanh). A *training set* is used to fit the parameters (weights). The network then “interpolates” to predict answers when confronted with data it has never seen before (the *test set*). For a general introduction to this field with emphasis on particle physics applications we refer the reader to, e.g., ref. [6].

When looking for new particles or resonances one almost always encounters the problem of computing invariant masses out of expected decay products. For

example, in the case of the intermediate vector boson W (produced in $p\bar{p}$ collisions) in its hadronic decay channel, $W \rightarrow q\bar{q} \rightarrow \text{hadrons}$, M_W is reconstructed relativistically from the momenta and assumed masses of the produced hadrons. A problem here is that the hadrons from the q and \bar{q} jets are not necessarily the only ones produced in the collision; there are also remnants from projectile hadrons and hadrons originating from bremsstrahlung. In addition there is of course the problem with noise in the detector.

In this letter we devise a method where a neural network learns to compute a mass from calorimetric information of hadronic decay particles. This kind of neural network application differs from the commonly used classification ones in the sense that the network learns to compute a real number (the mass) and not a set of binary decision values. A variation of the standard back-propagation learning algorithm [7] is used together with an appropriate choice of architecture.

2. The Monte Carlo data

The PYTHIA Monte Carlo [8] is used to generate

¹ E-mail addresses: lonnblad@desyvax (bitnet); lonnblad@apollo3.desy.de (internet).

² E-mail addresses: thepcap@seldc52 (bitnet); carsten@thep.lu.se (internet).

³ E-mail addresses: thepdr@seldc52 (bitnet); denni@thep.lu.se (internet).

data for the process

$$p\bar{p} \rightarrow W \rightarrow q\bar{q} \rightarrow \text{hadrons}, \quad (1)$$

where the hadrons emerge in a calorimeter detector with variables η , ϕ and E_{\perp} , strongly resembling the CERN SPS UA2 calorimeter with $|\eta| < 2$ and a resolution of 20 and 24 cells in the η - and ϕ -directions respectively.

The simulation includes the modeling of the underlying event [8] with a set of parameters according to the tuning found in ref. [9]. The selection of events is made by requiring two jets with $E_{\perp} > 10$ GeV in the $|\eta| < 1$ region.

3. Training set

To prevent the network from learning just one number (M_W), we generate W-bosons with fictive masses for training. These masses are flatly distributed with 300 events per GeV in the interval [40, 160] GeV, giving a total of 36×10^3 training examples.

4. Test sets

Two different data sets are used for testing; one with 10^4 events using the fictive flat mass distribution between 50 and 150 GeV, and one with a normal mixture of W- and Z^0 -events with $M_W = 80.0$ GeV and $M_{Z^0} = 91.2$ GeV. The former sample is used for monitoring the generalization performance during training and the latter to get a more realistic impression of how the network would perform on real data.

5. The neural network architecture and algorithm

We use a feed-forward architecture with one output node representing M_W , which is linear since it is not encoding a logical variable. Different numbers of hidden nodes and layers are explored. As expected, two or more hidden layers are needed to achieve peak performance. No further effort is made to optimize the hidden part of the architecture with respect to this task; our main goal is to provide a proof-of-concept for this novel application. Exploring different variants of input encoding and learning strategies for this

problem turns out to be illuminating, which is discussed below.

6. Input encoding

The following four different representations of the input calorimeter data are explored.

A: “Raw” calorimeter. A total of 480 input units, corresponding to the (η, ϕ) cells of the calorimeter, receive the E_{\perp} -values of each cell. These nodes are fully connected to a first hidden layer with $O(100)$ nodes, which are fully connected to a second hidden layer with $O(10)$ nodes, giving a total of $O(10^4 - 10^5)$ weights in the network.

B: “Raw” Calorimeter with receptive fields. As above but with the difference that each hidden node in the first hidden layer only covers a certain part of the calorimeter in connectivity (8×8 receptive fields). Each receptive field overlaps its neighboring fields and the total number of such fields is 312. Weights from the same relative position within the receptive fields are linked [6,7] such that they are updated identically, which is an effective way of encoding translational invariances. Each field is connected to three nodes in the first hidden layer, and the second hidden layer contains $O(10 - 10^2)$ nodes. The number of effective weights is thus decreased to $O(10^4)$.

C: Tower representation [3]. Take the E_{\perp} of the leading cell in the 20×24 matrix and assign it to the first node x_1 . Assign the η - and ϕ -coordinates to x_2 and x_3 respectively. Then take the second leading cell assign its E_{\perp} , η and ϕ to x_4 , x_5 and x_6 and so on for the first 30 calorimeter cells. This gives 90 input units. The hidden layers contained $O(10 - 10^2)$ and $O(10)$ nodes respectively giving $O(10^3 - 10^4)$ weights for the network.

D: “Intelligent” variables. The idea here is to use our physics knowledge of what may be useful variables for reconstructing the mass. The procedure we use is to first apply the LUCLUS [10] jet finding algorithm with parameters tuned such that there are always at least three jets being found. We then construct the invariant mass of the two largest jets (M_{12}), the three largest jets (M_{123}), and the four largest jets (M_{1234})^{#1}. The idea is that the invariant mass of the

^{#1} M_{1234} is set to 0 if less than four jets are found.

three largest jets is a good estimate of the W-mass in case a gluon has been radiated from one of the quarks. Furthermore, for each jet three variables are constructed, each of which is known to be sensitive to whether the jet stems from a quark or from a gluon [11]. These variables are

n_{90} : The minimum number of cells needed to account for 90% of the jets E_{\perp} . This is a real number; if eight cells contain 85% and nine cells contain 95%, n_{90} is set to 8.5.

n_c : The number of cells in the jet with $E_{\perp} > 1$ GeV.

s_c : The correlation of transverse energy deposited in neighboring cells in the jet, defined as

$$s_c = \sum_{i,j} \frac{|E_{\perp,i}^2 - E_{\perp,j}^2|}{(E_{\perp,i}^2 + E_{\perp,j}^2 + 1)}, \quad (2)$$

where i and j are nearest neighbors.

These variables are also constructed for the whole calorimeter. In addition the E_{\perp} of each jet as well as the total E_{\perp} and the total invariant mass of all towers found in the calorimeter is presented to the net, making a grand total of 24 input variables and $O(10^3)$ weights (see below).

In all cases the input variables have been rescaled to give an absolute value in the interval $[0,1]$, in order to avoid scaling problems in the weights. We also tried to use more than two hidden layers which in some cases speeded up the training but no significant improvement of the final performance of the net was found.

The variants with no preprocessing (A and B) give poor performance for our application, in contrast to the tower representation (C) and the "intelligent" variables encoding (D). The reason is that raw information requires many weights (parameters) to process, which with limited training data gives rise to poor generalization performance on the test set [12]. In other words one has a situation with too many parameters as compared to training examples (over-fitting). More quantitatively, using raw calorimeter data (A and B) yields generalization results that are comparable to what one obtains with the conventional method described above. In what follows we stick to alternative D above, which gives the best results.

7. Output encoding and learning strategies

We consistently use a summed square error measure for the back-propagation algorithm in our simulations.

$$E = \sum_p (o^{(p)} - t^{(p)})^2, \quad (3)$$

where $o^{(p)}$ is the output produced by the network for pattern p and $t^{(p)}$ the corresponding target value. It turns out that using the mass directly as target value ($t^{(p)} = M_t^{(p)}$) causes the network to systematically over-estimate the mass, since an error of 10% for a small mass is not as severe as it is for a large mass. What we really want to minimize is the width of the ratio between the reconstructed and the true mass, which can be achieved by instead using the logarithm of the true mass as target value for the network. The network then minimizes

$$\begin{aligned} E &= \sum_p [\log(M_o^{(p)}) - \log(M_t^{(p)})]^2 \\ &= \sum_p \left[\log\left(\frac{M_o^{(p)}}{M_t^{(p)}}\right) \right]^2, \end{aligned} \quad (4)$$

which for small errors approximately gives the width of the distribution of M_o/M_t :

$$E \approx \sum_p \left(1 - \frac{M_o^{(p)}}{M_t^{(p)}} \right)^2. \quad (5)$$

In both cases we normalize the target values to be between 0 and 1. In principle we could use a sigmoid response in the output layer. The idea would then be to start out with a low temperature (high gain), forcing the network to divide the patterns into high and low mass ones, then increasing the temperature gradually to fine-tune the net.

In the back-propagation algorithm the weights are changed by gradient descent. An alternative is the so-called *Manhattan* method [13] where the weights are changed by a fixed step size (the learning rate η) that is gradually decreased to 0 during training. The final result from the training is in principle independent of the strategy chosen but the latter has the advantage of being more tolerant to the choice of learning parameters than the former, especially when more than one hidden layer is used or when linear output nodes are used. In our calculations we have allowed the learning rate η to decrease geometrically,

$$\eta(t+1) = \eta(t)(1-\lambda), \quad (6)$$

where λ is a decay parameter.

8. The conventional approach

The “standard” procedure for reconstructing masses is to use a so-called cone algorithm to find two jets, defined as the two circular areas in the η, ϕ plane with radius $R^2 = \Delta\eta^2 + \Delta\phi^2$ with the largest E_{\perp} , and to calculate their invariant mass. In ref. [9] it is found that the optimum choice of R is 0.8. When we in the following compare with the ‘standard’ approach, we use the cone algorithm in JETSET [10] called LUCCELL. Although this is slightly different from the one used in ref. [9] it can safely be assumed that the choice of $R=0.8$ as the optimum value still holds.

9. Results

Based on the considerations above we use the following architecture and learning strategy

– *Inputs*: 24 input nodes according to coding strategy D above.

– *Hidden architecture*: Three hidden layers are used with 40 nodes in the first, 24 in the second and 10 in the third hidden layer.

– *Output*: One linear output node is trained to give

$0.5(\log M_W - 3.5) \in [0,1]$ for $M_W \in [40,160]$ GeV.

– *Initialization*: All weights are initialized randomly in $[-0.1,0.1]$.

– *Temperature*: The temperature for all nodes is set to 1.0.

– *Updating*: Patterns are selected randomly and the network is updated every 10th pattern using the Manhattan algorithm, starting out with $\eta=0.01$ and then lowering η with a factor 0.98 every epoch (= 36000 patterns) for 400 epochs.

The JETNET 2.0 package [14] is used for all the simulations. The performance of the network is monitored using the test set of 10^4 W-events with a flat mass distribution, measuring the width of the distribution in reconstructed mass divided by true mass. The final value of this width for the network is

$$\sqrt{\langle (M_W^{\text{net}}/M_W^0)^2 \rangle - \langle M_W^{\text{net}}/M_W^0 \rangle^2} = 0.15 \quad (7)$$

as compared with 0.19 when reconstructing the mass with the LUCCELL jet finding algorithm. In fig. 1 this distribution is shown for the test set with a normal mixture of events with W and Z distributed according to their true masses and widths. For LUCCELL the ratio has been rescaled so that $\langle M_{2\text{jet}}/M_W^0 \rangle = 1$. In fig. 2 we show the reconstructed W- and Z-mass distribution for the same test set.

The main reason why the neural network approach does better than the conventional method is that it

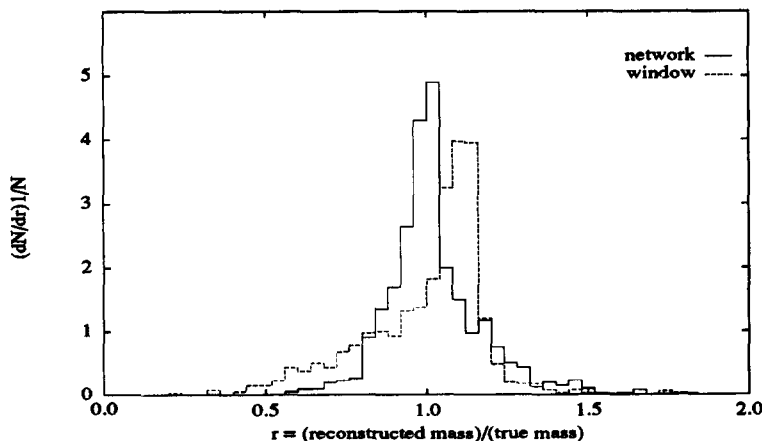


Fig. 1. The reconstructed mass ($M_{W,Z}$) divided by the true mass ($M_{W,Z}^0$) using the neural network method (full line) and the conventional “window” method (dashed line) with $R=0.8$.

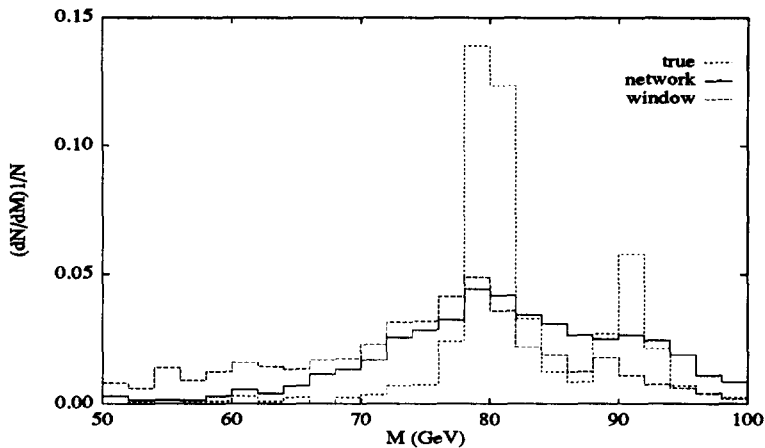


Fig. 2. The W - and Z^0 -mass peak for the true mass spectrum (dotted line), the spectrum reconstructed with the neural network (full line), and the spectrum reconstructed with the conventional "window" method (dashed line).

captures the gluon bremsstrahlung tails well and that it is more resistant to noise.

10. Summary

We have devised an efficient special purpose analog machine to reconstruct an invariant mass given calorimetric information. Being a neural network it is of "black box" nature and hence requires very little expertise knowledge to use and should be fairly fault tolerant and resistant to detector inefficiencies. The latter are easily included in the training phase.

The performance of the ANN algorithm is better than that of conventional ones; the reconstructed W -mass is more sharply peaked around the correct value and the resulting distribution is more symmetric, indicating that the algorithm handles gluon bremsstrahlung efficiently.

From experimenting with different architectures we strongly confirmed the fact that with limited sized training sets one needs to use preprocessed input data in order to avoid poor generalization (overfitting).

The method developed here could be of utmost importance when it comes to separating $H^0 \rightarrow W^+W^-$ produced at LHC/SSC energies from a background consisting of W +jets, $t\bar{t} \rightarrow W$ +jets and $b\bar{b} \rightarrow \ell$ +jets. In this case one needs to reconstruct M_w .

References

- [1] L. Lönnblad, C. Peterson and T. Rönngvaldsson, Phys.Rev. Lett. 65 (1990) 1321.
- [2] L. Lönnblad, C. Peterson and T. Rönngvaldsson, Nucl.Phys. B 349 (1991) 675.
- [3] P. Bhat, L. Lönnblad, K. Meier and K. Sugano, Using neural networks to identify jets in hadron-hadron collisions, preprint LU TP 90-13 in: Proc. 1990 DPF Summer Study on High energy physics research directions for the decade (Snowmass, CO, 1990), to appear.
- [4] L. Lönnblad, C. Peterson, H.Pi and T. Rönngvaldsson, Comput. Phys. Commun. 67 (1991) 193.
- [5] I. Scabai, F. Czakó and Z. Fodor, Quark and gluon jet separation using neural networks, ITP Budapest report 477 (1990).
- [6] C. Peterson and T. Rönngvaldsson, An introduction to artificial neural networks, preprint LU TP 91-23 (1991), in: Proc. Lectures 1991 CERN School of Computing (Ystad, Sweden), to be published.
- [7] D.E. Rumelhart, G.E. Hinton and R.J. Williams, in: Parallel distributed processing: explorations in the microstructure of cognition, Vol. 1, eds. D.E. Rumelhart and J.L. McClelland (MIT Press, Cambridge, MA, 1986).
- [8] H.-U. Bengtsson and T. Sjöstrand, Comput. Phys. Commun. 46 (1987) 43.
- [9] J. Alitti et al., Z. Phys. C 49 (1990) 17.
- [10] See e.g. B. Bambah et al., QCD generators for LEP, preprint CERN-TH.5466/89 (1989) [T. Sjöstrand, JETSET 7.3 program and manual]; T. Sjöstrand, Comput. Phys. Commun. 39 (1986) 347; T. Sjöstrand and M. Bengtsson, Comput. Phys. Commun. 43 (1987) 367.
- [11] J. Pumplin, Phys. Rev. D 44 (1991) 2025.

- [12] E.B. Baum and D. Haussler, *Neurulation Computation* 1 (1989) 151.
- [13] See e.g. C. Peterson and E. Hartman, *Neural Networks* 2 (1989) 475.
- [14] L. Lönnblad, C. Peterson and T. Rögvaldsson, *Pattern recognition in high energy physics with artificial neural networks - JETNET 2.0*, preprint LU TP 91-18 (1991), *Comput. Phys. Commun.* 70, No. 1 (1992) [Program and manual available via Email request].