

Report

A study of machine learning algorithms for reconstruction of missing energy in particle physics experiments.

Max Isacson, max.isacson@physics.uu.se
Mikael Mårtensson, mikael.martensson@physics.uu.se
Camila Rangel Smith, camila.rangel@physics.uu.se
Henrik Öhman, ohman@cern.ch

April 21, 2016

1 Introduction and Objective

High energy physics experiments have to deal with the so-called "missing energy", which refers to energy coming from particles that do not interact with the particle detector. These particles are usually weakly interacting neutrinos, but it could also come from physics beyond the Standard Model.

This missing energy is estimated by expecting conservation of energy and momentum in the system, and is reconstructed in the transverse plane. Although this quantity is generally attributed to undetected particles, apparent missing energy may also be caused by mis-measurement on the energy/momentum of the detected particles.

2 The Dataset

2.1 Structure

The dataset consists of simulated pp collision events, in which a charged Higgs is produced and decays as $H^\pm \rightarrow \tau\nu$. Each event is described by one set of observable variables and one set of unobservable variables.

Observable variables:

- $E_x^{\text{miss}}, E_y^{\text{miss}}$ — The x - and y -components of the missing energy.
- $P_{\tau_{\text{vis.}}}$ — The 4-momentum of the visible (hadronic) part of the τ decay.
- $P_{b_0}, P_{b_1}, P_{q_0}, P_{q_1}$ — The 4-momenta of the two b -jets and two light jets.

Unobservable variables:

- P_{ν_τ} — The 4-momentum of the neutrino from the charged Higgs decay.

- $P_{\bar{\nu}_\tau}$ — The 4-momentum of the neutrino from the τ decay.
- For simplicity, in this particular study it is assumed that the contribution from the neutrino from the τ decay to the total missing energy is negligible, and only focus on the neutrino from the charged Higgs decay as the unobserved variable.

2.2 Production

MG5_aMC@NLO [?] is used for the matrix element computation of $gg/q\bar{q} \rightarrow H^+$ and the event simulation. The events are then passed to PYTHIA8 [?] for the showering and hadronization, and for the $H^+ \rightarrow \tau\nu$ decay. Finally, the detector response is simulated using DELPHES [?] with an ATLAS-like geometry.

2.3 Estimated and true quantities

The true quantities for both the observable and the unobservable variables are available as output from the event simulation. The estimated values for the observable variables are reconstructed from the output from the detector response simulation.

3 Methods and Implementation

3.1 Bayesian Ridge Regression

A simple approach using linear regression techniques to predict the neutrino transverse moment was investigated. The algorithm with best performance is the Bayesian Ridge Regression, which introduces a zero mean spherical gaussian governed by a single precision parameter as prior probability distribution over the model parameters. The prior and noise parameter are estimated jointly during the fit to the data in several iterations. For an optimal result, 200 iterations are performed, it is found that a larger number of iterations does not improve the prediction.

The result presented here uses Python with Numpy [?], and Pandas [?] for general data processing, and scikit-learn [?] for the implementation of the algorithm.

3.2 Artificial Neural Network

An Artificial Neural Network (ANN) was used to predict the transverse momentum of the neutrino from the charged Higgs decay from the four-momenta of the detected decay products.

The ANN implemented here uses Python with SciPy [?], Numpy [?], and Pandas [?] for general data processing, and Keras [?] for the actual neural network. The number of input nodes is restricted to the 19 selected predictor variables and the output to the target dimension of 1. Since the network is used for regression, the activation functions between the last hidden layer and the output must be linear.

The data is scaled to avoid the tails of the ANN activation functions and improve the learning rate. The method used here is to scale it such that the minimum is -1 and

the maximum is 1 for each training set variable. The test set data is scaled using the scaling parameters computed from the training set. Another scaling method is to scale the mean to 0 and the standard deviation to 1, but this was found to slow down the training and give worst results.

Training is performed using stochastic gradient descent with the Nesterov method [?] and a mean-squared-error loss function. The learning rate was set to 0.1, the learning rate decay to 1×10^{-6} , and the momentum to 0.9.

A large number of configurations of the hidden layers were tested by varying the number of layers, the number of perceptrons in each layer, and the activation functions (sigmoid, tanh, and softmax). Sigmoid activation functions produced the best result for a fixed training period. A network with two hidden layers (i.e. 4 layers counting the 19 node input layer and 1 node output layer) with 25 and 10 perceptrons was found to be sufficiently complex. Using deeper networks, e.g. one with 4 hidden layer with 20, 35, 25, and 15 perceptrons, did not improve the result.

3.3 Support Vector Regression

Support vector machines can be used for regression, and the method is then called support vector regression (SVR). It employs two slack variables, $\xi_n \geq 0$ and $\hat{\xi}_n \geq 0$ and the corresponding error function to minimize is

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2$$

where C is the cost variable. The implementation used is the **SVR** class from **scikit-learn**, and it is an ε -SVR method. This implementation has four built-in kernels: RBF, sigmoid, linear, and polynomial. The selected predictors in the dataset are scaled to have a mean of 0 and a variance of 1 using the **StandardScaler** class from **scikit-learn**.

With the presented dataset, the selected SVR algorithm is not converging within a few minutes running time. The maximum number of iterations is therefore limited to 1000. A greater number of iterations is not found to yield significantly better results. The kernel that performs best is the RBF kernel. The value of the cost function is by default set to 1.0, but this setting yields a very poor resolution. Different values of the cost function are tested, and a value of $C = 100.0$ is found to be (near) optimal. For the soft margin the default value of $\varepsilon = 0.1$ is good.

Regardless of the choice of C and ε , the performance of the SVR is found to be no better than that of linear ridge regression. The resolution is not improved, and the shape of the true and reconstructed transverse component of the missing energy is not recovered by the prediction.

3.4 Trees and forests

Decision trees and ensemble methods, and in particular boosted decision trees (BDTs), have been a popular choice of method in particle physics. The method of decision trees are based on using multiple discriminating variables, and multiple choices in a binary

tree, to find the prediction that best describes the target variable. When using ensemble methods, such as extra trees, gradient boosting, or random forests multiple such trees are created, each with its own weight describing its "importance". The tree depth plays an important role when considering the trade-off between over-training and bias, and with the ensemble methods the number of trees is also an important parameter.

The performances of the classes `DecisionTreeRegressor`, `ExtraTreesRegressor`, `GradientBoostingRegressor`, and `RandomForestRegressor` in `scikit-learn` are examined. Different settings for the tree depth and (for the ensemble methods) number of trees are used. Near optimal values for the depth and the number of trees is found to be 15 and 100 for all methods and all ensemble methods respectively. The decision tree regression shows some improvement in resolution, but the shape of the transverse component of the missing energy is not recovered by the prediction. Neither in the ensemble methods is the shape of the missing energy recovered, but the resolution is improved. All three ensemble method show similar performance.

3.5 Gaussian process regression

Formally, a Gaussian process is a stochastic process defined as a probability distribution over functions $y(\mathbf{x})$, such that a set of function values y_1, \dots, y_N evaluated at inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ are jointly Gaussian distributed.

Given a training data set of targets $\mathbf{t} = (t_1, \dots, t_N)^T$ corresponding to inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$, regression can be performed on the new target t_{N+1} given new input \mathbf{x}_{N+1} by calculating the conditional probability $p(t_{N+1}|\mathbf{t})$. This will be a Gaussian distribution with mean and covariance given by

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (1)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \quad (2)$$

where the elements of \mathbf{C} are $C_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}$, the elements of \mathbf{k} are $k_n = k(\mathbf{x}_n, \mathbf{x}_{N+1})$, and $c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$. Here $k(\mathbf{x}, \mathbf{x}')$ is the kernel function, and β is the precision of the noise in the measurements, assuming that the noise ϵ_n for each measurement is isotropically distributed as $\mathcal{N}(\mathbf{0}|\beta^{-1}\mathbf{1})$.

The squared exponential was chosen as the kernel method, defined by

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|^2) \quad (3)$$

where the parameter θ is estimated using maximum likelihood, and is allowed to float between 10^{-4} and 10^{-1} . The precision of the noise was taken as $\beta = 100$.

Training and testing of the Gaussian processes was done using 3-fold cross validation on 8000 measurements. Due to the non-sparseness of Gaussian processes and memory constraints, more data than this was not feasible to use in the training.

4 Results

The metrics to measure the performance of the algorithms are the resolution and ratio of the predicted value with respect to its true value and the transverse mass of the tau

and neutrino system.

The transverse mass M_T of the τ and ν is defined as:

$$M_T = \sqrt{2p_T^\tau p_T^\nu (1 - \cos \Delta\phi_{\tau, \text{miss}})}, \quad (4)$$

where $\Delta\phi_{\tau, \text{miss}}$ is the azimuthal angle between the τ and the direction of the missing transverse momentum.

5 Discussion

This project focused on predicting the transverse momentum p_T of the neutrino from the charged Higgs decay by using the reconstructed four-momenta of the detected particles and the missing energy as inputs and training on the true p_T given by the Monte-Carlo simulation.

Out of all the Machine Learning (ML) methods tried, ANN and the bayesian ridge regression showed the best performance. The ANN shows the smallest bias in resolution with respect to the Bayes ridge regression but observes a significantly worse variance. All methods show a narrower transverse mass distribution, which is beneficial to discriminate signal over backgrounds.

The mass of the H^+ is the ultimate discriminant variable, but is not built due to the lack of longitudinal information in the detector, one option would be to target the H^+ mass directly and complement the missing information. This would require Monte-Carlo for a range of possible H^+ masses (3 mass-points were used in this study) and well for the background events.

As stated before, in this project it is assumed that the contribution from the neutrino from the τ decay to the total missing energy is negligible, and only focus on the neutrino from the charged Higgs decay as the unobserved variable. This approach is not correct for higher masses of the H^+ . A more complete study would be to target the complete three-vector of the missing energy or the (combined or separate) four-vectors of the two neutrinos. These options were briefly looked at but were considered to be out of the scope of this project. They will be pursued in more detail in future studies.