

Higgs Machine Learning Challenge with Neural Network Adversarial Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

The Higgs particle was theorized in 1964 and finally discovered in 2012 from the ATLAS and CMS experiments at CERN. This particle is the keystone of the Standard Model of particle physics, and provides the mechanism by which other particles acquire mass. Such discoveries are statistical in nature, and require vast amounts of information from the debris of particle collisions. Sifting and categorizing this data is a continual challenge, and in 2014 CERN and Kaggle invited the public to try their hand at such work, via the Higgs Machine Learning Challenge. A number of top entries, including the winner, used neural networks. An extended approach using the so-called "adversarial training" technique is presented herein.

1 Introduction

Particle physics arguably began with the discovery of the electron by J.J. Thomson in 1897 [?]. This began a century of subsequent rapid discoveries of new particles, next being the identification of the proton in 1919 [?] and then the neutron in 1932 [?], both by Ernest Rutherford. These three particles - electron, proton and neutron - are all that's needed to build the everyday chemical elements that comprise the apparent material universe. However, it was with early studies of cosmic rays that our understanding of fundamental particles began to change. One after the other more and more particles were discovered, with a broad spectrum of characteristics. Patterns slowly emerged, leading to the so-called Standard Model of particle physics, which encompasses not only particles of a material nature but also particles that transmit the very forces between them, thereby enabling their interactions.

As part of this framework, an altogether different particle was theorized to exist in order to explain why certain particles have mass. This came to be known as the Higgs particle, named after the author of the groundbreaking 1964 paper by one of the pioneers of this field [?]. So began a search that spanned nearly four decades, culminating with the discovery of the Higgs particle in 2012 from the Large Hadron Collider (LHC) at the CERN lab in Geneva, as measured by the ATLAS and CMS experiments [?]. The tell-tale signals were given by the physics properties of the Higgs particle's decay products, as the particle itself cannot be directly observed due to decaying so rapidly that it cannot travel out of its creation region and into the detector.

There are two main abilities of the experimental apparatus that allowed for the Higgs particle to be discovered, that were not had with earlier experiments. First and foremost is the particle collision energy, and second is the ability of the detector equipment to record and analyze data. Every particle has a characteristic mass energy, and this energy must be met by a device such as the LHC in order for the particle to be created. Energy is the raw ingredient of all particles, in other words, and different particles require different amounts of this ingredient. The Higgs particle happens to be a relatively massive particle compared to a number of others, and therefore requires a relatively large

amount of energy. A technological evolution over many years and many generations of device was necessary to develop the LHC, which is powerful enough to provide the Higgs creation energy.

Once created, a high mass particle such as the Higgs rapidly decays into other, lighter particles. It is these decay products that are detected by the machine surrounding the creation point, which are evaluated for such properties as momentum and electric charge, based on their trajectories through the machine, thereby leading to their classification. Originally, detectors were comprised of some continuous medium that would reveal a particle track that could be photographed and measured by hand. Things have since come a long way, with current detectors being formed of multiple layers of various electronics, each sampling the outgoing decay products in different ways. Despite the LHC providing an appropriate energy range for the creation of Higgs particles, their production is still exceedingly rare compared with the production of a number of other particles. An exceedingly large number of creation events is therefore required to compensate for the rarity, and the particle tracks simply cannot be assessed individually as used to be possible in photographic detectors from previous generations. The vast stream of electronic readout must instead be dealt with computationally and statistically, and a great deal of time and effort is spent by scientists determining how best to seek the signals of interest.

Machine learning has inevitably become an invaluable tool in this task, which, following the reconstruction of particle properties from the raw signal, is simply one of classification, namely it is a question of identifying signal versus background. In the case of the Higgs particle, its various decay channels are not unique, and are in fact shared by the decays of other heavy particles. Seeing a certain set of decay products does not then give indisputable proof of one particular parent particle or the other. But there are nonetheless subtle signs within the data when viewed *en masse*, which machine learning algorithms can be trained to recognize.

In 2014, CERN scientists collaborated with the Kaggle machine learning organization, developing an open competition for the public to try and develop a machine learning classifier to distinguish the Higgs particle from other heavy particles in the tau-tau decay channel (explained further in 2) [?]. This was the Higgs Machine Learning Challenge (HMLC), which at the time of its release was the most popular Kaggle challenge up to that point, with a total of 35,772 uploaded entries submitted by 1,785 teams (1,942 people) comprised of participants from around the world, some working within the physics community, others simply interested in computing [?]. The winning entry came from a sole programmer with no advanced physics training, who developed a neural network model that was clearly superior to anything else [?].

Many technical fields are developing incredibly fast at present, and machine learning is no exception. Even just more than one year on from the Higgs machine learning challenge, new ideas and techniques have become available. This project seeks to explore one such technique that has started to generate a lot of interest of late, namely that of adversarial training, and to apply it to the HMLC.

2 Higgs Machine Learning Challenge

Despite machine learning being a key component to many physics analyses at CERN, there is relatively little development of these techniques within the organisation itself [?]. When one technique is found to work and is well understood, and, crucially, it is well understood how to interpret the resulting physics results, then that technique gets entrenched and remains in play for a long time. This is not a bad thing in many ways, as by having staple workhorse codes it is possible for physicists to spend a greater fraction of their time analysing the data in search of physics.

However, machine learning is a field in its own right, and has its own experts and development trends quite separate from any of its applications, in science or otherwise. The last decade in particular has seen an enormous rise in both the ability of machine learning and the applications [?]. An ever increasing fraction of the developed world's activities has some computational aspect, and with that comes more data and more need to analyse that data for patterns. Coupled with this is the increase in the ability of hardware to run ever more complicated algorithms, making theoretical work from previous generations of computer science suddenly more relevant and useful.

The Higgs Machine Learning Challenge (HMLC) was developed in order to tap into the knowledge of this field straight from the enthusiasts, in order to port that knowledge into the particle physics domain. This is not only beneficial for CERN, but also the machine learning community at large,

which demands as many real world problems as possible in order to develop the most relevant algorithms [?]. Sourcing good data is a job in itself, and if there's one thing that CERN is good at it's generating data. The HMLC was therefore a mutually beneficial arrangement for the data generation experts and the data processing experts.

The HMLC has two datasets - one fully labelled dataset for training, and another unlabelled dataset for testing. The training dataset has information about 250,000 labelled events, and the testing dataset corresponds to 550,000 unlabelled events. There are only two labels involved, namely "signal" and "background". In a particle physics collider, two input particles (protons in the case of the LHC for ATLAS and CMS experiments) are accelerated and collided at high speed, resulting in a high energy state that can form any particle which "costs" that amount of energy or less. Call this created particle the *initial state*. This initial state often has a number of possible decay channels to some set of *final state* particles, and some initial states share similar final states. In the case of the Higgs particle, one decay channel is that to a pair of lighter tau particles, with these then decaying to yet lighter particles still. These events are those labelled as "signal". Such final states are possible from the decay of other initial states such as W and Z particles (which are the carriers of the weak force) and top particles (which are the heaviest of the quark family) [?]. These events are those labelled as "background".

It is the subtle differences in the ensemble physics variables that we seek to distinguish through machine learning. Both training and testing datasets are defined by a feature set comprised of 30 different physics variables, representing quantities such as the mass, energy and momentum of decay products that entered the detector or were otherwise inferred [?, ?]. These entities are either (1) electrons, (2) heavier variants of electrons (the muon and tau particles), (3) jets of particles comprised of quarks (which are particles within the family of the constituents of protons and neutrons), and (4) so-called missing energy, representing very light and weakly interacting particles called neutrinos, which are not directly detected but inferred via a deficit in the energy accounting [?].

Physics details aside, this format of problem is ubiquitous and well known to machine learning, and as such requires no knowledge or understanding of what any of the features actually are. The HMLC dataset is actually a vastly simplified version of the data actually encountered by CERN physicists, who have to carefully reconstruct such high level output from a very low level input. The HMLC dataset was carefully tailored to be simply enough for public use, but still complicated and realistic enough to promote machine learning techniques that could be of actual use [?]. Furthermore, the dataset is in fact the product of a very fine Monte Carlo simulation, of the sort that is actually used for understanding the background in real data. The signal to background ratio in the HMLC dataset is much larger than that in real data, which requires many millions of events in order to provide statistically significant findings.

3 HMLC Winner and Initial Impulse

As mentioned in §1, there were many hundreds of entrants to the original Kaggle competition, with the overall winner being a single programmer with no professional physics training, specializing instead in computational development and consulting [?]. In fact, the physicist with the best score came in eighth place overall, highlighting the disparity in the computational skill sets between the professions, as should be somewhat expected. It was, after all, with such an outcome in mind that the challenge was run in the first place.

The winning entry was based on the neural network approach. The final iteration of this entry utilized 3 layers of 600 nodes each, finished off with 2 softmax output units (one for signal, one for background) [?]. Some feature preprocessing was performed on the input data, normalizing each feature to have zero mean and unit standard deviation, with some also being log-transformed. Further to the given features, four more were derived from some of the given angular features, and five more were derived from some of the given mass features, so giving nine extra hybrid quantities [?]. The algorithm was trained by minimising cross-entropy, with a correctness probability then assigned to each classification, and with only those above a certain threshold retaining their determined class labels at any one iteration [?]. Further to this, a range of model variants were employed simultaneously via *bagging*, with each variant given a different random subset of the training data and with their outputs then being averaged.

All these elements of the winning entry were included for good reasons. The effects of each were carefully examined over a four month development period, along with many others that were found to be of no benefit and discarded. The speed with which computer science is evolving means that even now, little more than a year after the close of the original Kaggle competition, there are new things to try that were not part of the winning formula. Soon after the start of this current project, the Kaggle winner was contacted for advice based on their past experience and their understanding of the machine learning field in general. They proposed that a so-called generative approach could be applicable to this task, and so initial research was then undertaken in this direction, which became the seed of this project's central concept.

4 Adversarial Training

The generative approach originally suggested by the Kaggle winner regarded Generative Adversarial Networks (GANs), which have been the subject of much interest over the last year or so, following work based on image recognition published in 2014 by Goodfellow *et al.*, based at the University of Montreal [?]. The idea is to train not only one machine learning system, but two. One system tries to discriminate the classes of a given set of training data, and the other system tries to generate new data examples that can fool the former system into passing them as real. The discriminator D is positively reinforced for every example it correctly classifies, and for every fake example it deems as such, with the generator G being positively reinforced for every example it makes which passes off as real data. In practice this can be done by setting up D and G as multilayer perceptrons, using backpropagation and dropout training techniques with piecewise linear units, with the reinforcement criteria set up as a competitive "minimax" game between the two systems, hence the "adversarial" in the GAN acronym.

Although the original work by Goodfellow *et al.* was intended to create an adept generator, for the particular case of creating lifelike images, the fact that both systems must become stronger simultaneously, if either are to do so at all, means that the discriminator also becomes as adept as possible. It is this latter fact that is of relevance to this work, in which we require the best possible discriminator for the HMLC. The GAN process, insofar as interests in the discriminator go, essentially comes down to creating as difficult a situation as possible for the discriminator, thereby making it more robust to examples outside of the range of the real training data. In other words, if D can be made to work well for simulated worst-case scenarios, then it should perform well on any typical or atypical real world examples it may encounter when used in practice.

Following basic research into the methods and codes used by the Montreal group, further direct advice was sought from the group itself. The short timeframe of the project described herein meant that rapid evaluation of feasibility was necessary, and the GAN approach, although interesting and promoted by the Kaggle winner themselves, seemed to present a challenge in itself to come anywhere near putting into practice as needed. Thankfully, the authors of the original GAN paper responded swiftly with further information, and a much simplified adversarial approach was suggested.

The full GAN concept described above is still useful to understand the adversarial process in general, namely having one part of a system increase the challenge for another part that we wish to strengthen. Neither opponent in the adversarial system is aware of the other, as such, but rather each component simply affects the other's error minimisation, effectively supplying an additional regularization term. It is this fact that allows for simplification, for instead of developing a full minimax-playing two-component system in order to obtain a better discriminator, as in the present HMLC case, we can instead begin with some discriminator and then modify its error function to emulate the regularization effect of a sophisticated generator.

In practice this is still done by use of training examples that are not in the supplied training data, just not in the same way as the GAN approach. There exists a dedicated generative code in a GAN system that takes only random values as input, with a discriminator that is rewarded for spotting the produced fakes as well as possible, as described. The simpler approach generates new examples by taking the actual training data and perturbing it slightly, ideally on the order of the data noise level, if known, which is much simpler than generating convincing examples from random values. The perturbed data should, ideally, be classified in exactly the same way as the unperturbed data i.e. it should not be different enough from the unperturbed data to warrant any other classification. This is an economical way of generating high quality adversarial examples, of the type that an actual

generator aims to produce from scratch, namely within the data noise limit and, therefore, hopefully indistinguishable from the real thing.

A simple neural network based on a multilayer perceptron works well for this procedure. In this case the error function presents itself as an update to the network weights with every training example, with the examples assessed one by one for a pre-defined number of epochs. The weights are randomly initialized via a normal distribution with zero mean. Feeding a training example through the network leads to a set of errors at the output layer, which are backpropagated to obtain adjustment derivatives for each layer of weights. These derivatives should reduce the overall error when applied to the existing weights. This is the standard neural network backpropagation procedure.

However, in the case of simplified adversarial training the weight update does not actually occur at this stage in the algorithm. The error is instead backpropagated once more, giving adjustment derivatives for the input layer itself. The input layer is not being reiterated towards some goal, as is the case with the weights, but rather the input layer derivatives are the means by which a training example is perturbed to produce an adversarial example. It is, in fact, not the value of this perturbation but rather its *sign*, multiplied by some constant, that is used for the actual perturbation applied to the real training example. This is referred to as the "fast gradient sign method". The constant involved is ideally on the order of the data noise level, or at least small enough to not perturb the training example too much.

The adversarial example is then fed forward through the same network (i.e. with the same weights) as used for the real training example. Once again the output errors are calculated and backpropagated, giving a second set of adjustment derivatives for each layer of weights. What we then have are two sets of derivatives for each layer of weights, one set from the real training example and another set from the perturbed adversarial example. Each set of derivatives say by how much the network must be corrected to cater for the example in question. If the network treats both the training example and its adversarial counterpart the same, which is the intention, then the two sets of derivatives should be the same. If this is not the case then the two sets of derivatives will differ. By updating the weights by the average of the two sets of derivatives, the corrected network should not only become more adept in general, but also become more likely to treat further examples and their adversarial partners more equally.

The above can be mathematically summarized via the following modified cost function, in which J is the original cost function, Θ is the model parameters, \mathbf{x} is the input vector, y is the target label, α is the mixing parameter (set to 0.5 for equal treatment of real and perturbed training examples), and ϵ is the perturbation constant (ideally on the order of the data noise, as mentioned):

$$\tilde{J}(\Theta, \mathbf{x}, y) = \alpha J(\Theta, \mathbf{x}, y) + (1 - \alpha) J(\Theta, \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\Theta, \mathbf{x}, y)), y)$$

5 Baseline Algorithms

In order to observe the effects of the chosen adversarial training approach, and to observe the evolution of the implementation, a number of baselines were created. First was the adaptation of two existing classifier codes, one being driven by a Support Vector Machine (SVM) algorithm and the other by a Kernelized Perceptron (KP) algorithm. The former operates as part of the Python SciKit-Learn framework, and the latter was in fact based on the MatLab code developed as part of SFU CMPT419/726 assignment 2, but translated into Python in order to be run outside of the MatLab environment. This was soon found to be necessary for the sake of running on significant fractions of the given training data, which is large enough to give memory management issues when running on personal machines, as so with the license-restricted MatLab suite. Given the relative simplicity of the HMLC dataset, there was nothing to suggest that these relatively simple SVM and KP approaches would be invalid, and indeed they performed as desired.

5.1 SVM

As SciKit-Learn has a mature SVM implementation in `sklearn.svm.SVC`, it was a clear choice to use with the HMLC dataset. The SVC class supports several kernel functions by default, each with its own set of parameters. To find optimal parameters, a 5-fold cross-validation was performed using `sklearn.grid_search`. To facilitate this parameter tuning, grid search was run on a 50,000 example

subset of the entire dataset. The following table lists the kernels and parameters tried during the grid search. Note that γ and C are ranges. For C the range is exponents, while the gamma values are spaced evenly on a log scale, starting and ending with the values in the table.

kernel	C	γ	coef0	degree
radial-basis function	10 – 100,000	0.00001 – 100,000	N/A	N/A
sigmoid	10 – 100,000	0.001 – 10,000	0, 1	N/A
polynomial	10 – 100,000	0.001 – 10,000	0,1	1, 2, 3

Having chosen an optimal kernel and parameters, the entire training set of 250,000 examples was randomly split into 225,000 training samples, and 25,000 validation samples. The SVM was trained using a radial-basis function (RBF) kernel, with parameters $C = 1,000$ and $\gamma = 0.00024883895413685354$. The overall score obtained was 73.744%. For the signal class, precision was 83 percent, and recall was 30 percent. For the background class, precision was 72%, while recall was 97%.

5.2 Kernelized Perceptron

The kernelized perceptron algorithm has the ability to evaluate a range of kernel parameters through cross-validation in order to find the optimum for a given kernel function. However, in order to have a direct comparison with the SVM, it was chosen to use the RBF kernel with the γ -parameter value found to be optimum for this kernel with the SVM (the C -parameter mentioned above is specific to the SciKit-Learn implementation of the RBF, and so not relevant for the MatLab kernelized perceptron). A scan of the gradient descent step size η was performed, first with un-normalized data and then with normalized data. Three such runs were done for each value in the range $10^{-5} \leq \eta \leq 10^{-1}$, but due to time restrictions only one such run was done for each value in the range $10^0 \leq \eta \leq 10^5$, and the results are shown in Fig1. Only 10% of the full 250,000 example HMLC training data was used here, with 20,000 examples used for training and 5,000 examples used for testing. All runs were done on a single node of the SFU High Energy Physics (HEP) group computing cluster.

It can be seen that there is similar performance for both the training and testing data when normalized, and although performance is generally improved in both cases for un-normalized data, the improvement is more significant for the training set than the testing set. These observations indicate that normalization is not only a hindrance for this algorithm, but that it is somehow affecting the expected performance, which should typically be better for training data compared to testing data. Furthermore, there seems to be a greater spread in the performance values of normalized data, showing lower stability with respect to the stochastic procedure. Although not conclusive due to relatively low statistics, it could be that there is a peaked trend in the performance for normalized data, where the average seems to tend to a maximum around $\eta = 1$. This may also be true of the un-normalized data, but the effect is much more subtle and would require more work to assess properly. The best performance on un-normalized testing data was 76.98% for $\eta = 10^{-5}$, and the best performance on normalized testing data was a slightly larger 77.54% for $\eta = 1$, but, for the reasons discussed, this latter value is not as impressive in light of the general negative aspects of normalization seen here.

6 Main Algorithms

With the intention of developing a neural network with adversarial training in mind, a logical first step was to modify an existing neural network already in hand. Once again the SFU CMPT419/726 assignments came in useful here, with the third and final assignment involving the completion of a simple MatLab neural network for handwritten digit recognition. This code was adapted for use with the HMLC dataset, and given some further refinements such as optional training data normalization. Satisfied that the code worked as intended in this state, the adversarial components were then introduced. These additions sit entirely within the inner calculation loop, following the existing error backpropagation code lines but before the weight update step, as described in §4. By simply setting the perturbation size ϵ to zero, the adversarial calculations are rendered null and void, and the code performs the same as prior to their inclusion, as should be so.

Due to time restrictions, this MatLab code was unfortunately not translated into a more general language such as Python, which, as mentioned in §5, would have allowed ease of use of a greater

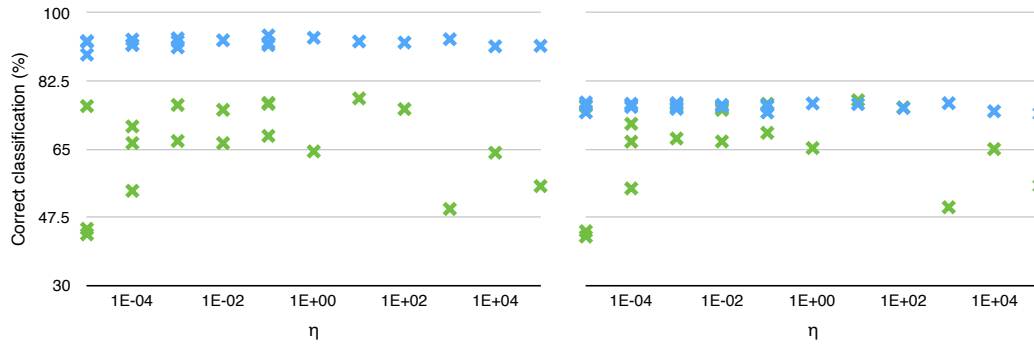


Figure 1: Left: Training data; Right: Testing data
Blue: Un-normalized; Green: Normalized

fraction of training data. However, despite the input limitations, the MatLab neural network still performed rather well.

!!! MatLab NN details !!!

The MatLab neural network was intended to be no more than an advanced baseline scenario, and to provide a development platform for the adversarial training components that could then be incorporated into something more sophisticated overall. To this end, following an in-depth survey of available tools, the Python-based Theano framework was chosen as a basis for the final stage of this work. Theano is a specialized package for performing tensor computations efficiently by treating them symbolically. As such, it is very useful in machine learning applications. It is a competent package which is well supported and, importantly, well documented, and was a natural choice following the progression of work up to this point. Further, as the modification to the cost function J required access to several gradients of the cost function, Theano's symbolic integration allowed for quick experimentation with many different cost functions. An example Theano neural network code was taken as a basis and adapted for the HMLC dataset. Specifically, a Multilayer Perceptron was written using Theano to classify the HMLC dataset. The network consisted of 30 node input layer, normalized, a hidden layer with 600 nodes and tanh as an activation function, and an output layer with two nodes, using softmax to read the output. The network was trained on a 2013 Macbook Pro with 2 cores and 8GB of ram, for 100 epochs. As with the baseline algorithms, the data was split into 225,000 training samples and 25,000 validation samples. Two cost functions were used in testing. The results are summarized below:

cost	non-adversarial	adversarial
Negative log likelihood	73.368	74.552
Cross entropy	65.5	65.5

7 Conclusion

Contributions

- All: Survey of available tools; General research; Poster editing; Paper editing
- Darren Temple: MatLab/Python kernelized perceptron adaptation; MatLab neural network adaptation; Initial adversarial training implementation; Communication with external contacts
- Oren Shklarsky: Python SVM adaptation; Python/Theano neural network adaptation; Advanced adversarial training implementation
- Morgan Jenkins: Interpretation of research material; Development of adversarial training algorithm; Assistance with Python/Theano neural network
- Lola the dog: Running; Jumping; Licking heads; General distractions

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

References