

Searching for the Higgs Boson Particle using Data Analytics

Ananya Choudhury^{*}
Dept. of Computer Science
Virginia Tech
Blacksburg, U.S.A
ananya@vt.edu

Vivek Bharath Akupatni[†]
Dept. of Computer Science
Virginia Tech
Blacksburg, U.S.A
vivekb88@vt.edu

Vignesh Adhinarayanan[‡]
Dept. of Computer Science
Virginia Tech
Blacksburg, U.S.A
avignesh@vt.edu

ABSTRACT

Abstract goes here

1. INTRODUCTION

Introduction goes here

2. BACKGROUND

Physics background

3. METHODOLOGY

3.1 Data Preprocessing

In some situations, it becomes difficult to measure the physical properties of a particle such as momentum and energy accurately. In our original dataset, as many as 177000 out of 250000 instances had missing attributes. We considered three approaches to deal with missing values. First, we tried ignoring instances with missing attributes. But this resulted in very few useful instances. Second, we tried to replace the missing values with the mean/median. However, this biases the experiments. Finally, we decided to adopt a method known as multiple imputations [?] which replaces the missing values with a random number that follows the distribution for that attribute. We use the *Amelia* [9] package to perform this task.

3.2 Feature Engineering and Selection

The raw dataset includes 17 features. From these 17 basic features, 13 additional features were derived. These 13 derived features describe some property of the particle and requires knowledge of physics. The features are provided by physics. Their description is given in the appendix. From these 30 features, a subset is selected based on the following

^{*}Ananya worked on.

[†]Vivek worked on.

[‡]Vignesh worked on.

factors. First, we decide the ability of the feature/attribute to distinguish between signal and background. The distribution of different attributes is shown in Fig. 1 for signal and background separately. Second, we try to avoid using features correlating with each other while building the classifier. Fig. 2 shows the correlation matrix of the features. The final set of features selected differs for each classifier. The details are given in their respective subsections.

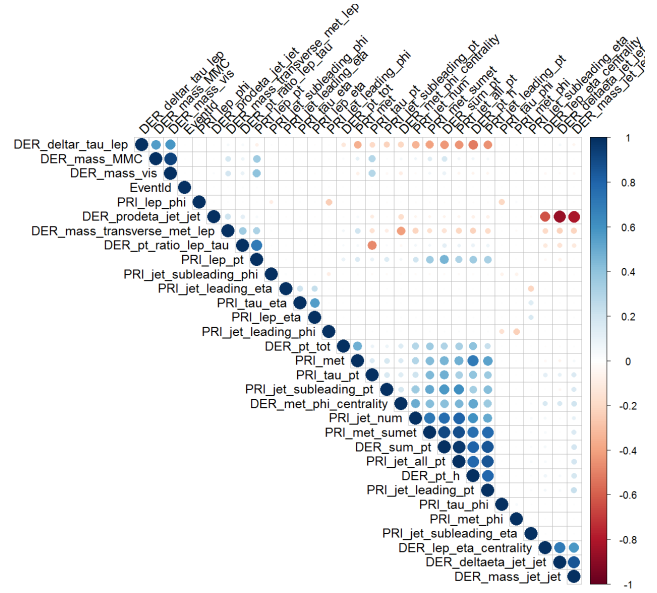


Figure 2: Correlation matrix. Dark blue indicates features that shows strong positive correlation. Dark red indicates features that show strong negative correlation.

3.3 Data Analytics Techniques

In this section, we describe the various classification schemes we explored. For each technique, we also describe the parameter settings explored in brief.

3.3.1 Bayesian Classifiers

Naive Bayes. This classifier is based on the Naive Bayes technique developed by John et al. [10]. In this technique <describe technique here>. <Describe parameter tuning here>.

3.3.2 Functions-based Classifiers

Logistic Regression. This classifier is based on the Ridge estimation technique developed by Cessie et al. [6]. In this technique <describe technique here>. <Describe parameter tuning here>.

Linear Discriminant Analysis

Quadratic Discriminant Analysis

3.3.3 Tree-based Classifiers

Simple CART. This classifier is based on the classification and regression trees (CART) technique developed by Brieman et al. [5]. In this technique <describe technique here>. <Describe parameter tuning here>.

3.3.4 Instance-based Classifiers

k-Nearest Neighbor. This classifier is based on the k-nearest neighbor (kNN) technique by aha et al. [1]. In this technique <describe technique here>. <Describe parameter tuning here>.

3.3.5 Deep Learning

3.3.6 Meta Classifiers and Ensemble Methods

Rotation Forest. Based on Rotation Forest technique developed by Rodriguez et al. [15].

Decision Stump tree with ADA Boosting. Based on ADA Boosting developed by Freund and Schapire [8]. Used in conjunction with Decision Stump Tree classifier.

Decision Stump tree with Multi-Boosting. Based on the MultiBoosting technique developed by Webb [17].

REP Tree with Bagging. Based on the bagging technique developed by Brieman [4].

Classification via Clustering and Regression. We simply cluster the raw dataset and mark certain clusters as signal and others as noise. Prediction based on the distance of the new data point to the centroid of the two cluster groups.

4. RESULTS

4.1 Results

Bayes Classifiers

Rule-based Classifiers

Tree-based Classifiers

Instance-based Classifiers

Deep Learning

Ensemble Methods

4.2 Discussion

5. RELATED WORK

Detecting exotic particles in high-energy physics (HEP) using data analytics techniques instead of the traditionally used physical detectors [14] is not new. Cutts et al. are among the first to use neural networks to identify interesting events in HEP experiments [7]. This was quickly followed by several attempts in improving the classification accuracy of neural networks [11, 13]. Apart from the widely popular techniques based on neural networks, only decision trees were explored for a long time. Bowser-Chao and Dzialo used binary decision trees to detect top quarks and compared their results with neural networks [3]. The conventional wisdom was that neural networks were by far the best when it comes to classification in HEP until Roe et al. came along and projected boosted decision trees as an alternative to artificial neural networks [16]. By combining several *weak* classifiers, Roe et al. showed that it is possible to obtain better accuracy than a neural network. However, more recently, Baldi et al. showed that a deep learning neural network with several hidden layers outperforms the boosted decision tree [2]. While a number of classification techniques have been developed and applied over the last several years, the HEP community has so far explored only neural networks and boosted decision trees in any depth. This led to the development of statistical packages for the HEP community such as StatPatternRecognition so that several other techniques based on Rotation Forest, Discriminant Analysis etc. could be explored [12]. Despite this effort, no documented work exists in HEP where alternative techniques are explored even though other techniques are thought to be inferior.

Since Baldi et al. [2] work is closest to ours, we describe it in detail here. The authors in this paper used deep learning methods of neural network to find exotic particles in high energy particle colliders. Deep learning models are neural networks with multiple hidden networks. Current techniques like shallow models which are single hidden layer feedforward network fail to capture all features. The deep learning model here is used on 2.6 million training examples and 100,000 validation examples. The model is a five-layer neural network with 300 hidden units in each layer, learning rate of 0.05, and a weight decay coefficient of 0.00001. Testing is done on 500,000 examples. For Higgs benchmark, Area Under the Curve (AUC) - complete for deep neural network is 0.88 and for shallow neural network = 0.81.

6. CONCLUSION

Conclusion goes here

7. REFERENCES

- [1] D. W. Aha, D. Kibler, and M. K. Albert. Instance-based learning algorithms. *Mach. Learn.*, 6(1):37–66, Jan. 1991.
- [2] P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5:4308, July 2014.

- [3] D. Bowser-Chao and D. L. Dzialo. Comparison of the use of binary decision trees and neural networks in top-quark detection. , 47:1900–1905, Mar. 1993.
- [4] L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [5] L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA, 1984. new edition [?].
- [6] S. L. Cessie and J. C. V. Houwelingen. Ridge estimators in logistic regression. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 41(1):pp. 191–201, 1992.
- [7] D. Cutts, J. Hoftun, A. Sornborger, R. Astur, C. R. Johnson, and R. T. Zeller. The use of neural networks in the d0 data acquisition system. *Nuclear Science, IEEE Transactions on*, 36(5):1490–1493, Oct 1989.
- [8] Y. Freund and R. E. Schapire. Experiments with a new boosting algorithm. In *Thirteenth International Conference on Machine Learning*, pages 148–156, San Francisco, 1996. Morgan Kaufmann.
- [9] J. Honaker, G. King, and M. Blackwell. Amelia II: A program for missing data. *Journal of Statistical Software*, 45(7):1–47, 2011.
- [10] G. H. John and P. Langley. Estimating continuous distributions in bayesian classifiers. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, UAI’95, pages 338–345, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
- [11] L. Lönnblad, C. Peterson, and T. Rönqvist. Finding gluon jets with a neural trigger. *Phys. Rev. Lett.*, 65:1321–1324, Sep 1990.
- [12] I. Narsky. StatPatternRecognition: A C++ Package for Statistical Analysis of High Energy Physics Data. *ArXiv Physics e-prints*, July 2005.
- [13] C. Peterson. Track finding with neural networks. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 279(3):537 – 545, 1989.
- [14] J. Pinfold. Searching for exotic particles at the {LHC} with dedicated detectors. *Nuclear Physics B - Proceedings Supplements*, 78(1–3):52 – 57, 1999. Advanced Technology and Particle Physics.
- [15] J. J. Rodriguez, L. I. Kuncheva, and C. J. Alonso. Rotation forest: A new classifier ensemble method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(10):1619–1630, 2006.
- [16] B. P. Roe, H.-J. Yang, J. Zhu, Y. Liu, I. Stancu, and G. McGregor. Boosted decision trees as an alternative to artificial neural networks for particle identification. *Nuclear Instruments and Methods in Physics Research A*, 543:577–584, May 2005.
- [17] G. I. Webb. Multiboosting: A technique for combining boosting and wagging. *Machine Learning*, Vol.40(No.2), 2000.

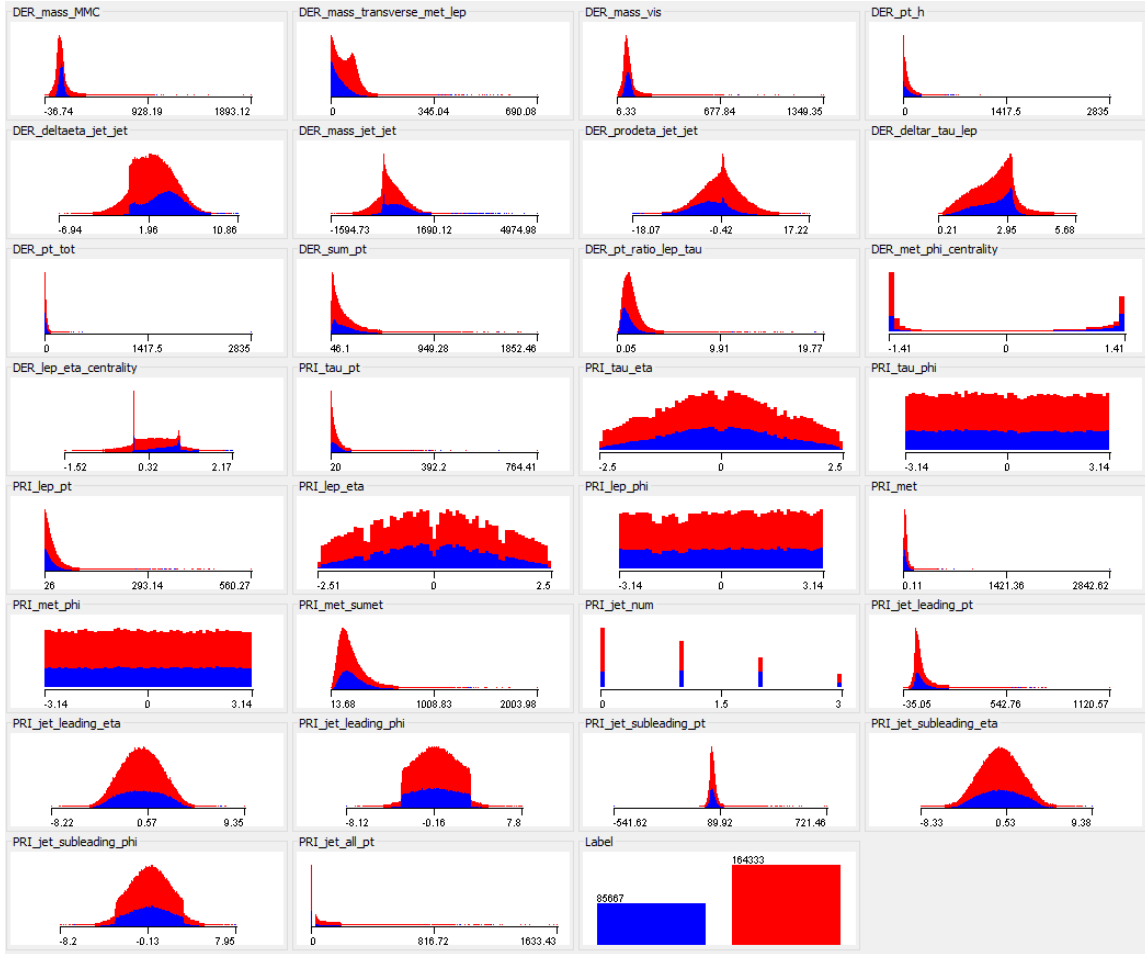


Figure 1: Distribution of different attributes for signal. Blue represents signal and red represents background.