

Snowmass2021 - Letter of Interest

Detecting New Physics as Novelty

Thematic Areas:

- (TF07) Collider phenomenology
- (CompF3) Machine Learning

Authors:

Xuhui Jiang (Department of Physics, The Hong Kong University of Science and Technology, Hong Kong S.A.R.) [xjiangaj@connect.ust.hk]

Aurelio Juste (Institut de Física d'Altes Energies, Edifici Cn, Facultat de Ciències, Universitat Autònoma de Barcelona, Spain, Institució Catalana de Recerca i Estudis Avançats, Spain) [Aurelio.Juste.Rozas@cern.ch]

Ying-Ying Li (Theoretical Physics Department, Fermi National Accelerator Laboratory, U.S.A.) [yingying@fnal.gov]

Tao Liu (Department of Physics, The Hong Kong University of Science and Technology, Hong Kong S.A.R.) [taoliu@ust.hk]

Abstract: Novelty (anomaly) detection is a task of machine learning to detect novel events without a prior knowledge. The deep-neural-network-based techniques for novelty detection recently received high attention and have been proposed for searching for unexpected signals of new physics at colliders. As a foot stone of many such techniques, the evaluators of data novelty can be roughly classified into two classes, namely isolation-based and clustering-based, depending on whether the evaluation of each testing event correlates with those for the others. In this study, we demonstrate that a complementarity generically exists between these two classes of evaluators. We are now developing a new class of novelty evaluators which allow us to synergize the pros of these two classes of evaluators. For illustration, we plan to apply this new design to analyzing the $t\bar{t}\gamma\gamma$ data at LHC, with the $t\bar{t}h \rightarrow t\bar{t}\gamma\gamma$ Higgs production and the gravity-mediated supersymmetry with a final state of $t\bar{t}\gamma\gamma$ serving as novel events in this context.

After the triumph of the Higgs discovery^{1;2}, null findings in NP have motivated the design of new analysis strategies which allow the unexpected NP to be detected in a more model-independent way and with a broader coverage in theory space, and hence complement the current ones extensively used at LHC.

In the science of ML, this involves a well-known task - novelty (or anomaly) detection (for a review, see, e.g.,³): detect novel events without a prior knowledge. This implies that there is no data of the signal pattern available for model training. Various algorithms have been applied for the NP search⁴⁻¹⁰ in the recent years,. Many of these novelty evaluators can be approximately classified into two classes⁴, according to the emphasis of evaluation.

- Isolation-based. The novelty response of a given testing point is evaluated according to its distance to or isolation from the distribution of known-pattern data in the feature space. All of the other testing points are irrelevant in this process.
- Clustering-based. The novelty response of a given testing point is evaluated according to the clustering around this point on top of the distribution of known-pattern data in the feature space. Some other testing points potentially in the same cluster of unknown-pattern data are relevant in this process.

Here the distribution of the known-pattern data in both cases can be figured out by taking either Monte-Carlo simulation (semi-supervised ML) or data extrapolation (fully unsupervised ML).

In our previous work⁴, \mathcal{O}_{iso} and \mathcal{O}_{clu} , representing each of both classes, are defined with the method of k -nearest neighbors. In short, \mathcal{O}_{iso} measures the difference between the mean distance of a testing data to its k nearest neighbors(kNN) and the average of the mean distances defined for its k nearest neighbors(kNN's kNN), while \mathcal{O}_{clu} will evaluate novelty response of the testing point by comparing its local densities in the training and testing datasets. Furthermore, The constructions for the isolation-based and clustering-based evaluators are not unique. For instance, the reconstruction error of autoencoder^{5;6;8} in essence is isolation-based, while the likelihood ratio used in^{7;9;10} shares similar spirit to that of \mathcal{O}_{clu} .

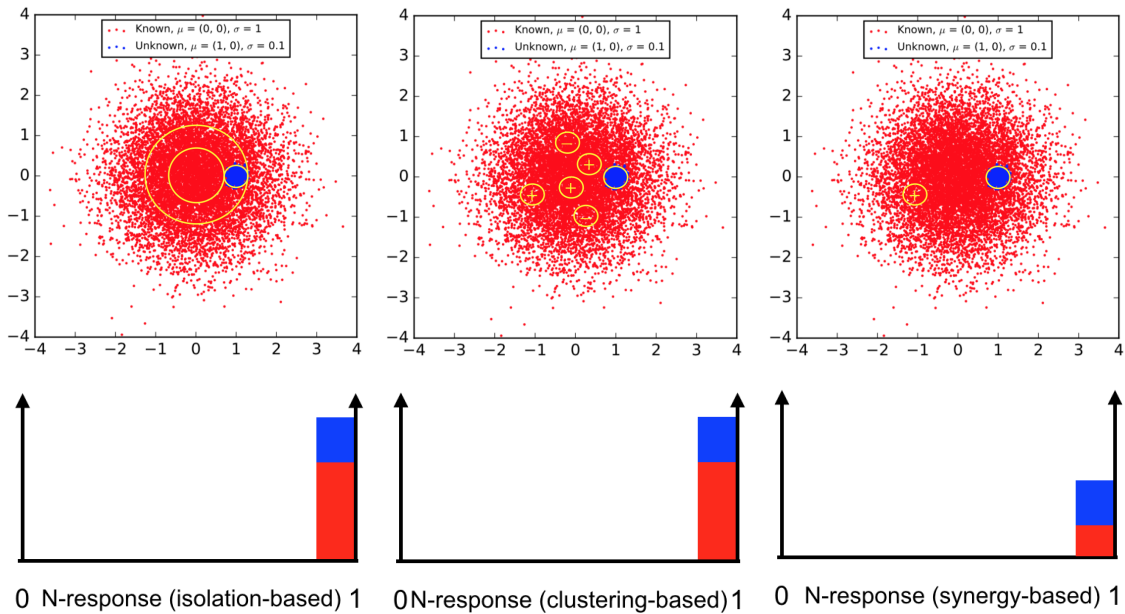


Figure 1: Novelty response of data to \mathcal{O}_{iso} , \mathcal{O}_{clu} and \mathcal{O}_{syn} - cartoon demonstration. Here red and blue points are of the known and unknown patterns, respectively.¹¹

Yet, to improve their effectiveness and efficiency in a general sense, it is important to synergize these two classes of novelty evaluators. The reason is demonstrated in Fig. 1, with 2D Gaussian samples. Here the red and blue points represent the known-pattern and unknown-pattern data, respectively, with their novelty response being evaluated by the isolation-based and clustering-based evaluators. In this figure, the signal bin of both (as is shown in the bottom panels) receives contributions of red points from not only the true signal region (blue area in the top panels), but also some non-signal region, at the 2D plane. For the isolation-based evaluators, e.g., \mathcal{O}_{iso} , the relevant non-signal region is the ring area between the two yellow circles (excluded the original signal (blue) region) in the top panels. For the clustering-based ones, the relevant non-signal region is characterized by upward fluctuations. If these two classes of the evaluators are well-synergized, one would expect that only the intersection between the two non-signal regions may contribute to the signal bin significantly, resulting in an improved S/B . To this purpose, we design the third class of novelty evaluators:

- Synergy-based. The novelty response of a given testing point is evaluated according to the synergization of its novelty responses or the novelty responses of the datasets to both isolation-based and clustering-based novelty evaluators.

In ⁴, the synergy-based evaluator is naively defined as:

$$\mathcal{O}_{\text{syn}} = \sqrt{\mathcal{O}_{\text{iso}}\mathcal{O}_{\text{clu}}}, \quad (1)$$

i.e., the geometric mean of \mathcal{O}_{iso} and \mathcal{O}_{clu} .

However, \mathcal{O}_{syn} has its own disadvantages. There is no reason to regard the geometric mean as the optimal solution to synergizing the 2 classes of evaluators. Besides, \mathcal{O}_{clu} breaks the statistical independence of the testing points and further the Gaussian/Poisson properties of the dataset, since it evaluates novelty response of a given testing point based on its correlation with some other testing points in the dataset. While calculating the significance based on the novelty response to \mathcal{O}_{clu} and hence \mathcal{O}_{syn} , one needs to extract out the one-sigma definition in significance, using, e.g., a p -value method.

Thus we are motivated to invent a more generic synergy-based novelty evaluator $\mathcal{O}'_{\text{syn}}$. A binary supervised neural network is design to synergize \mathcal{O}_{iso} and \mathcal{O}_{clu} automatically where its output is defined as $\mathcal{O}'_{\text{syn}}$. Though its definition is based on the novelty responses of testing data to \mathcal{O}_{iso} and \mathcal{O}_{clu} , as a non-linear function, $\mathcal{O}'_{\text{syn}}$ evaluates the novelty response of each testing point independently. Therefore, it preserves the Gaussian/Poisson statistics of the testing dataset successfully. More than that, different from manually selecting the geometric mean as the novelty measure, we expect such a design of synergy-based evaluator would open a broad avenue for us to optimize the performance of novelty detection at colliders.

We initiate with two-dimensional Gaussian samples, as a proof of concept, to mimic signals of different kinematic patterns. $\mathcal{O}'_{\text{syn}}$, as a new design, will be studied as the integration of the isolated-based and clustering-based evaluator. We expect to see a complementarity between and then will generalize the algorithm to physical cases. As a concrete application, we propose to analyze the $tt\gamma\gamma$ data at LHC. The $tth \rightarrow tt\gamma\gamma$ Higgs production and the gravity-mediated supersymmetry with a final state of $tt\gamma\gamma$ will serve as novel events, which represent the signal patterns with a sharp resonance and a broad distribution of $m_{\gamma\gamma}$, respectively.

References

- [1] G. Aad *et al.*, “Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC,” *Phys. Lett.*, vol. B716, pp. 1–29, 2012.
- [2] S. Chatrchyan *et al.*, “Observation of a New Boson at a Mass of 125 GeV with the CMS Experiment at the LHC,” *Phys. Lett.*, vol. B716, pp. 30–61, 2012.
- [3] M. A. Pimentel, D. A. Clifton, L. Clifton, and L. Tarassenko, “A review of novelty detection,” *Signal Processing*, vol. 99, pp. 215–249, January 2014.
- [4] J. Hajer, Y.-Y. Li, T. Liu, and H. Wang, “Novelty Detection Meets Collider Physics,” 2018.
- [5] M. Farina, Y. Nakai, and D. Shih, “Searching for New Physics with Deep Autoencoders,” 2018.
- [6] T. S. Roy and A. H. Vijay, “A robust anomaly finder based on autoencoder,” 2019.
- [7] A. De Simone and T. Jacques, “Guiding New Physics Searches with Unsupervised Learning,” *Eur. Phys. J.*, vol. C79, no. 4, p. 289, 2019.
- [8] T. Heimel, G. Kasieczka, T. Plehn, and J. M. Thompson, “QCD or What?,” *SciPost Phys.*, vol. 6, no. 3, p. 030, 2019.
- [9] R. T. D’Agnolo, G. Grosso, M. Pierini, A. Wulzer, and M. Zanetti, “Learning Multivariate New Physics,” 12 2019.
- [10] B. Nachman and D. Shih, “Anomaly Detection with Density Estimation,” *Phys. Rev. D*, vol. 101, p. 075042, 2020.
- [11] X. Jiang, R. Aurelio.Juste, Y.-Y. Li, and T. Liu, “Detecting New Physics as Novelty,” (ongoing).