

# MLAnalysis: An open-source program for high energy physics analyses

Yu-Chen Guo, Fan Feng, An Di, Shi-Qi Lu, Ji-Chong Yang\*

*Department of Physics, Liaoning Normal University, Dalian 116029, China*

*Center for Theoretical and Experimental High Energy Physics, Liaoning Normal University, Dalian 116029, China*

---

## Abstract

We present a python-based program for phenomenological investigations in particle physics using machine learning algorithms, called **MLAnalysis**. The program is able to convert LHE and LHCO files generated by **MadGraph5\_aMC@NLO** into data sets for machine learning algorithms, which can analyze the information of the events. At present, it contains three machine learning (ML) algorithms: isolation forest (IF) algorithm, nested isolation forest (NIF) algorithm, kmeans anomaly detection (KMAD), and some basic functionality to analyze the kinematic features of a data set. Users can use this program to improve the efficiency of searching for new physics signals.

Source code: <https://github.com/NBAlexis/MLAnalysis>

## Program summary

Program Title: MLAnalysis

Programming language: Python3.8 and above

Nature of problem: With the continuous accumulation of experimental data, the research of high energy physics needs to process a large amount of data. ML methods can help us to improve the effect and efficiency of data analysis. Converting the data from experiments or Monte Carlo (MC) simulated events into data sets available for ML has become an important requirement. A program platform is needed for data preparation, as well as the application of various ML algorithms to improve the selection capability of target events and the efficiency of particle identification.

Solution method: Supply an event analysis platform that supports ML approaches. The program is able to convert LHE and LHCO files into data sets that can be used for ML algorithms, and apply data preparation. In the data preparation step, the program transforms the raw data into a format that can be used to train and test machine learning algorithms, optimizes the adaptabilities and generalization capabilities of algorithms. The program offers several algorithms, including IF, NIF, and KMAD, which provide NP model independent and standard model effective field theory operator independent methods to optimize event selection strategies.

**Keywords:** Particle Physics Phenomenology, Analysis, Recasting, Machine Learning

---

## 1. Introduction

The search for new physics (NP) beyond the Standard Model (SM) is one of the most important tasks of high energy physics (HEP). In most cases, due to the good agreement between experimental measurements and the SM predictions, NP signals are expected to be rare events and their kinematic behaviors are different from that of the SM. So it is necessary to optimize event selection strategies (ESSs) with the help of the kinematic characteristics. However, in some cases, it is difficult to efficiently suppress the backgrounds by kinematic cuts, so a better method for selecting signal events is needed.

Due to the fact that the NP signals are usually rare, a large number of events need to be analyzed both for experimental data and Monte Carlo (MC) simulation cases. Machine learning (ML) methods can be helpful to improve the effect and efficiency of data analysis, and have been used in various aspects of HEP [1–13]. With continuous improvements, ML has played an important role in particle identification [2–4, 13], searching for NP signals [14–22] and studying the polarization of final particles [23–26], etc. The need for a program framework that can make ML algorithms more accessible is motivated.

In a typical ML algorithm, there are several main procedures involved, including: data preparation, model selection, training, evaluation, hyperparameter tuning, and deployment. In these steps, data preparation is a critical step in the ML process which is often overlooked but can have a significant impact on the performance of the model. If we want to use ML algorithms to analyze experimental data or MC simulation events, we need to pre-process these data through data preparation. Data preparation is the transformation of raw data into a form that is more suitable for modeling. Data preparation involves data cleaning, data transformation, feature engineering. Data cleaning is to remove any missing values, correcting errors or inconsistencies. Data transforms are used to change the type of distribution of data variables and normalizing the data to fall within a certain range. Feature engineering is the process of selecting, transforming, and creating features that will be used to train the ML algorithms. Features are the variables that represent the input data and can be used to make predictions about the output. The goal of feature engineering is to create features that are relevant, informative, and non-redundant. This might involve selecting the most important features using techniques such as correlation analysis or feature importance scores, transforming the features using techniques such as scaling, normalization, or one-hot encoding, or creating new features that capture interactions or higher-order relationships between the original features. Data preparation is a critical step in a ML process, that can significantly impact the performance and accuracy of the final model.

Recently, with the advancement of the study of the processes at the hadron collider with multi-neutrinos in the final states [21, 22] and the optimization scheme of event identification [15, 20], we integrate the programs used in these studies into automatic program tools, called `MLAnalysis`. The `MLAnalysis` allows users to efficiently perform pre-defined and custom analyses of event files generated by the MC event generators such as `MadGraph5_aMC@NL0`. The code has been tested in python 3.8 and above.

The importance of `MLAnalysis` lies in its ability to help researchers unlock the value of the experimental or MC data. Many researchers collect large amounts of data in the

---

\*Corresponding author

Email addresses: `ycguo@lnnu.edu.cn` (Yu-Chen Guo), `yangjichong@lnnu.edu.cn` (Ji-Chong Yang)

course of their work, but may not have the expertise or resources to analyze it using ML. This program can help to bridge this gap by providing a user-friendly tool for data transformation and feature engineering that can help researchers to extract insights and knowledge from their data. In addition to its data transformation and feature engineering features, this program also includes basic functionality for data processing and several ML algorithms. This opens unlimited possibilities concerning the level of complexity which can be reached, being only limited by the programming skills and the originalities of the users.

The structure of the paper is as follows. Section 2 overviews the structure and functions of `MLAnalysis` package. Then section 3 introduces the ML algorithms provided. Section 4 shows some general usages by examples. Our summary can be found in section 5.

## 2. Program overview

The `MLAnalysis` acts as a bridge between phenomenology studies and ML algorithms by transforming the data from experiments or MC into a format that can be recognized by ML algorithms, and also assists with feature engineering. Some of the key features of the program include data cleaning and preprocessing, data transformation and normalization, feature selection and extraction, and data visualization. These features can help to automate and streamline the process of preparing the data for ML, saving time and reducing the risk of errors. `MLAnalysis` also includes basic functionality for data processing and several ML algorithms. By including these additional features, `MLAnalysis` provides a more complete suite of tools for working with phenomenological studies and ML algorithms.

The purpose of this program is to make it easier for researchers to work with experimental data or the data from MC in the context of ML. By providing a user-friendly interface and a suite of tools for data transformation and feature engineering, `MLAnalysis` can help to remove some of the barriers to entry for those new to the field of ML, and to make ML more accessible to a wider range of researchers in HEP.

### 2.1. Data Structure

The data structure of `MLAnalysis` mainly includes five classes. The `EventSet` is the uppermost container, which has only a list of objects with `EventSample` as the class. The `EventSet` represents a set of collision events. The `EventSample` is a class lists all particles of a collision event, this data structure contains a list of objects with `Particle` as the class. The class `Particle` contains information about a particle, including particle type, four-momentum, mass, Particle Data Group identifier (PDG-id), etc. The four-momentum is represented by the class `LorentzVector`. The class of `Martix4x4` is used for operations on a `LorentzVector`, such as rotation and Lorentz boost.

In Table 1, we show the functions of these classes and the types of input parameters of functions, where “float”, “str” and “int” stand for the floating-point number, string, and integer, respectively.

Table 1: The functions of the classes of **MLAnalysis** and the types of input parameters of these functions.

		LorentzVector	
		values: list	
		Azimuth(): float	
		Et(): float	
		MakeWithRapidity(pseudoRapidity: float, azimuthal: float, pt: float, mass: float)	
		Mass(): float	
		P3d()	
		PseudoRapidity(): float	
		Theta()	
		V3d()	
		Y(): float	
		Particle	
		PGDId : int	
		index : int	
		momentum	
		particleType	
		status	
		DebugPrint(sep: str): str	
		SetLHCOOtherInfo(nTrack: float, bTag: float, hadEm: float)	
		Matrix4x4	
		values : NoneType, list	
		MakeBoost(v3velocity)	
		MakeOne()	
		MakeRotation(degree: float, x: float, y: float, z: float)	
		MakeRotationFromTo(v3from, v3to)	
		MakeRotationFromToV4(v4from: LorentzVector, v4to: LorentzVector)	
		MakeZero()	
		MultiplyMatrix(otherMatrix)	
		MultiplyVector(vector: LorentzVector): LorentzVector	
		EventSet	
		events : list	
		AddEvent(event: EventSample)	
		AddEventSet(eventSet)	
		DebugPrint(i: int)	
		GetCopy()	
		GetEventCount(): int	
		EventSample	
		particles: list	
		AddParticle(particle: Particle)	
		DebugPrint(): str	

## 2.2. Import a data set

Performing a phenomenological analysis on the results provided by MC generators or by experiments always starts with the reading of a set of event samples. Currently, the types of event files supported by **MLAnalysis** include the Les Houches Event (LHE) files at the parton-level, and LHC Olympics data (LHCO) files at the reconstruction level. The files “LesHouchesEvent.py” and “LHCOlympics.py” in the “Interfaces” folder are used for importing. We classify all particle types into seven categories: “jet”, “electron”, “muon”, “tau”, “photon”, “intermediate” and “missing”. It should be noted that, in an LHE file the “missing” refers to each neutrino, whereas in an LHCO file it is the sum of transverse momenta of all neutrinos.

```
def LoadLHCOlympics(fileName: str) -> EventSet:
```

This function is fed with a full path to the LHCO file, and will return an object with **EventSet** as the class.

```
def SaveToLHCO(fileName: str, event: EventSet, realLHCO: bool = True):
```

This function is fed with a full path to the LHCO file and an **EventSet**. The **EventSet** is read from an LHE or an LHCO file or built from the code. There might be incoming particles and intermediate particles. Besides, there might be multiple neutrinos in an **EventSample**. When “realLHCO” is turned on, only the outgoing particles are saved, and the neutrinos are combined into a missing transverse momentum.

```
def LoadLesHouchesEvent(fileName: str) -> EventSet:
```

This function is fed with a full path to the LHE file, and will return an object with **EventSet** as the class.

```
def LoadLargeLesHouchesEvent(fileName: str, debugCount: bool) -> EventSet:
```

Similar to **LoadLesHouchesEvent**, but reads an LHE file line after line. When “debugCount” is turned on, a message will be printed after each **EventSample** is loaded.

The support for other file formats is on the way.

## 2.3. Data cleaning and cuts

Data cleaning is one of the steps in data preparation. Besides, the ESS based on traditional kinematic cuts can be used to extract signals from backgrounds. In both cases,

a set of tools for cuts is necessary. The files in the “CutAndExport” folder present the cut mechanism and some cuts for common usages. To implement a cut, the “CutEvents” function can be used.

```
def CutEvents(eventSet: EventSet, cutFunction):
```

The “CutEvents” function is fed with an `EventSet`, and an object “cutFunction”. It is required that the interface “Cut” is implemented for the “cutFunction”.

```
def Cut(self, eventSample: EventSample) -> bool:
```

The format of the interface “Cut”. The “Cut” function should be fed with an `EventSample`. Whether this `EventSample` should be cut off depends on the returned value of the “Cut” function. If the “Cut” function returns “True”, this event sample should be cut off, otherwise, should remain in the `EventSet`.

As an example, a typical class for “cutFunction” is presented as follows.

```
1 class PhotonNumberCut:
2     """
3     If cutType = 0, cut all with photons > parameters[0]
4     If cutType = 1, cut all with photons < parameters[0]
5     If cutType = 2, cut all with photons not in parameters
6     """
7
8     def __init__(self, cutType: int, parameters):
9         self.cutType = cutType
10        self.parameters = parameters
11
12    def Cut(self, eventSample: EventSample) -> bool:
13        photonCount = 0
14        for particle in eventSample.particles:
15            if 0 == particle.particleType:
16                photonCount += 1
17        if 0 == self.cutType:
18            return photonCount > self.parameters[0]
19        if 1 == self.cutType:
20            return photonCount < self.parameters[0]
21        return photonCount not in self.parameters
```

The code defines a class named “PhotonNumberCut”, which has an initializer method that takes two arguments: “cutType” and “parameters”.

The “cutType” argument is an integer that determines how the event sample should be filtered based on the photon count, and “parameters” is a list of values that is used by “cutType”.

The “Cut” method takes an “EventSample” object as an argument and returns a boolean value. This method counts the number of particles in the event sample that have a “particleType” of 0 (which represents photons) and compares it to the “parameters” based on the “cutType”. If “cutType” is 0, it returns true if the photon count is greater than the first element of “parameters”. If “cutType” is 1, it returns true if the photon count is less than the first element of “parameters”. If “cutType” is 2, it returns true if the photon count is not in the list of “parameters”.

The “CutAndExport” folder also contains some unique cuts proposed in our anomalous quartic gauge coupling (aQGC) and the neutral triple gauge coupling (nTGC) studies [19–22, 27–31]. The “Applications” folder contains the complete project files for these studies.

Sometimes, the generic features might however not be sufficient according to the needs of the users. The projects in the “Applications” folder can provide references for users, and their analysis can be realized by using and modifying the relevant codes. Taking advantage of the Python interface, a user can define his own physics analysis in an efficient, flexible and straightforward way.

### 3. Methods of Machine Learning Algorithms

The NP signals are generally rare and kinematically different compared with the SM background. Thus, the search for NP can be considered as anomaly detection (AD), then, ML algorithms can be used to find kinematically anomalous events. `MLAnalysis` contains ML algorithms for AD, and therefore can be used in the search of NP.

#### 3.1. Isolation Forest Algorithm

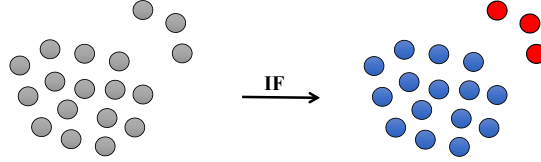


Figure 1: The IF algorithm can identify anomalous data. The signal points are the anomalous points colored in red, and the background points are colored in blue.

IF algorithm [32] is an unsupervised ML algorithm with linear complexity to deal with AD problems, which can effectively deal with large-scale multi-dimensional data. It is good at finding the data which are “few and different”. Compared with normal samples, such anomalous samples are more easily isolated. IF can be used to search for anomalous samples by constructing a binary tree structure, where the anomalous samples are closer to the root node compared with the normal samples (see Fig. 1). Such binary trees are called Isolation Trees (iTrees).

Here we use the vector boson scattering (VBS) process  $pp \rightarrow jj\ell^+\ell^-\nu\bar{\nu}$  with leptonic decay at the 13 TeV LHC as an example, which has been studied using IF [19]. The  $t\bar{t}$  production with b-jet mistagged is also considered as the background. The signal and background events are generated by `MadGraph5_aMC@NLO` [33], with a parton shower by `Pythia82` [34] and a detector simulation by `Delphes` [35]. IF algorithm is applied to calculate the anomaly scores for signal and background events. It proves that the IF algorithm has the ability to distinguish between the SM signal and the background of the  $t\bar{t}$  process, as shown in Fig. 2.

When using the IF algorithm, there is no need to know what kind of NP signal the data set contains, and there is no need to optimize the parameters according to the characteristics of the NP signal. In other words, it is possible to select NP signals without a priori knowledge of the NP models.

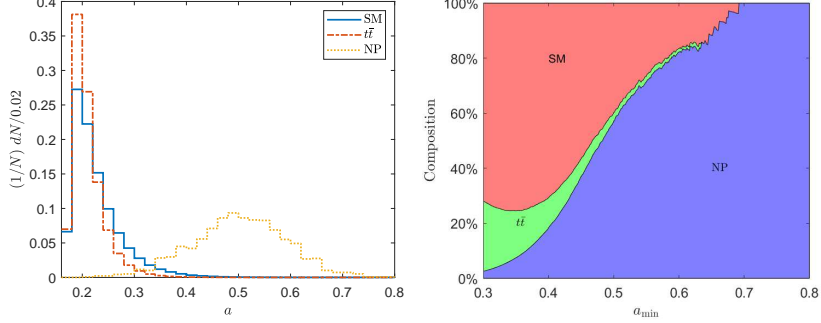


Figure 2: Normalized distributions of anomaly scores denoted as  $a$  (left) and compositions of the selected events with  $a > a_{\min}$ , where  $a_{\min}$  is a threshold to select the events (right).

### 3.2. Nested Isolation Forest Algorithm

Although the IF algorithm could select NP signals, AD-based algorithms were no longer applicable when the interference between NP and the SM is important. To illustrate this, we use two dimensional points to describe this problem in Fig. 3. As can be seen from Fig. 3, the signal in this case will increase the density of the data point distribution rather than been more easily isolated. According to the IF, the anomaly scores would be higher in the area with low density. Thus, in the case of Fig. 3, the anomaly scores of some points would be reduced when NP presents. If the distribution of the anomaly scores for the SM events is taken as a reference benchmark, the variation of density in phase space can be measured by the change of anomaly scores. Then the existence of an NP signal can be detected in this way.

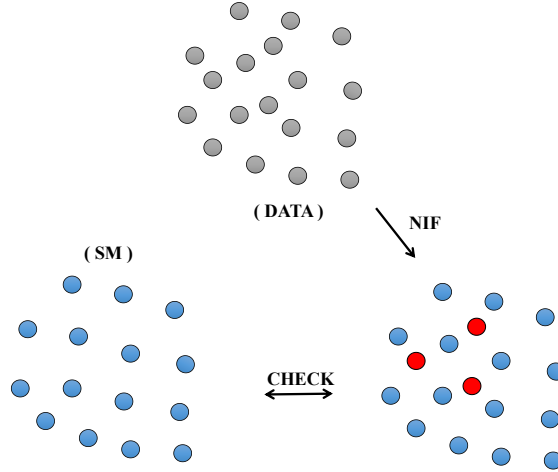


Figure 3: The distribution of signal points overlaps with the background, which is no longer an AD problem. By comparing the data distribution of the SM, NIF can select the points that change the distribution.

Based on this idea, we propose an unsupervised ML algorithm, which is called nested

anomaly detection. When an IF algorithm is nested, it is then nested IF (NIF). First, the MC simulation data set of the SM is used as the training data set, which is marked as  $S_{\text{SM}}$ , and the anomaly scores  $a_{\text{SM}}$  of each event are obtained by the IF. Then, the anomaly scores  $a_{\text{data}}$  for each event in the target data set  $S_{\text{data}}$  are obtained by the IF. Finally, the closest events in the phase space between the target data set and the training data set are matched. The change in anomaly score for each pair of events is  $\Delta a^i = a_{\text{data}}^i - a_{\text{SM}}^i$ . Here the distance is defined as  $d = \sqrt{\sum_{ij} (p_j^i - q_j^i)^2}$ , where  $p$  and  $q$  are the 4-momenta of the particles in  $S_{\text{data}}$  and  $S_{\text{SM}}$ , respectively, and  $p_j^i$  and  $q_j^i$  are the  $i$ -th component of the 4-momentum of the final-state particle  $j$ .  $\Delta a$  is the indicator used to detect the existence of NP signals. When the NP signal exists in the events, this indicator will be less than 0. We can adjust the sensitivity of the NIF algorithm to the NP signal by setting the maximum value of  $\Delta a$ .

NIF algorithm not only inherits the advantages of IF, which is independent of NP models and SM effective field theory (SMEFT) operators, but also solves the problem that AD can not be handled. It has an intelligible operation mechanism and almost no adjustable parameters. In addition, the NIF program framework of `MLAnalysis` can be used not only for IF but also for any algorithm that can quantitatively measure the abnormal degree of each event, which results in a good generality. The search of nTGCs in the process  $e^+e^- \rightarrow Z\gamma$  at the  $e^+e^-$  colliders is an example that the NIF algorithm works well [20].

### 3.3. K-means anomaly detection method

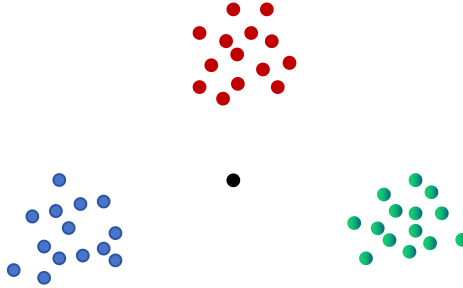


Figure 4: The K-means algorithm can divide the background events (the colored points) into several groups. The anomalous events are expected to be far away from all centroids of the groups. An example of the anomaly is depicted as the black point.

Another algorithm along similar lines to IF is an anomaly detection algorithm based on K-means. K-means-based anomaly detection (KMAD) is a method of anomaly detection that uses the K-means clustering algorithm to identify anomalous points in a data set.

The K-means algorithm is a popular clustering algorithm that partitions a data set into K clusters, where each point belongs to the cluster with the closest centroid [36]. In KMAD, the K-means algorithm is used to cluster the normal points in the data set, and the distance from a point to the nearest centroid of the clusters is used to quantify the degree of anomaly of the point [37].

The basic idea of KMAD is that anomalous points are those that are far from all of the centroids of the normal clusters. The distance of a point from the nearest centroid can be used as a measure of its degree of anomaly. That is, the background is categorized using the K-means algorithm, such that the distribution of the background events can be described and sampled by a set of centroids around which the events are distributed. If the anomalous events are kinematically different from the ones of the background, the anomalous events will be farther away from all the centroids. Therefore, if a point is farther away from all the centroids than a certain threshold, it can be considered an anomaly. An illustration of KMAD is shown in Fig. 4.

In Ref. [37], the KMAD is used to study the NP contribution in the process  $\mu^+\mu^- \rightarrow \gamma\gamma\gamma$ , which is an example that the KMAD works well.

#### 4. Typical usage

To summarize, data preparation and feature engineering are critical steps in ML that cannot be overlooked. The quality and relevance of the features can have a significant impact on the performance of the model. Therefore, it is important to spend sufficient time and effort in preparing the data and engineering the features to ensure the best possible outcome.

In this section, we present a few examples on how to use the `MLAnalysis` in the data preparation phase, as well as how `MLAnalysis` can be used directly to select event signals.

##### 4.1. Build a data set for KMAD, IF and NIF algorithm

Data preparation and data culling are necessary aspects before ML algorithms. In the case of the process  $\mu^+\mu^- \rightarrow \gamma\gamma\gamma$ , for example, the final state sometimes contains less than two photons due to detector simulations. At this point, frequently, retaining only events containing three final-state photons improves the accuracy of signal screening. For the same reason, the final state sometimes contains more than three photons, at which point we can select the hardest three photons and compose their four-momenta into a 12-dimensional vector for each event. The 12-dimensional vectors are then to be stored in a “csv” file for use in the next KMAD, IF or NIF algorithms. The following code is an example of using `MLAnalysis` to accomplish the above.

```

1  def ChooseEventWithStrategy(allEvents: EventSet, count: int, tag: int):
2      result = []
3      idx = 0
4      while len(result) < count:
5          theEvent = allEvents.events[idx]
6          largestPhotonIndex = -1
7          largestPhotonEnergy = 0
8          secondPhotonIndex = -1
9          secondPhotonEnergy = 0
10         thirdPhotonIndex = -1
11         thirdPhotonEnergy = 0

```

```

12     for theParticle in theEvent.particles:
13         if theParticle.particleType == ParticleType.Photon:
14             PhotonEnergy = theParticle.momentum.Momentum()
15             if PhotonEnergy > largestPhotonEnergy:
16                 thirdPhotonIndex = secondPhotonIndex
17                 thirdPhotonEnergy = secondPhotonEnergy
18                 secondPhotonIndex = largestPhotonIndex
19                 secondPhotonEnergy = largestPhotonEnergy
20                 largestPhotonIndex = theParticle.index - 1
21                 largestPhotonEnergy = PhotonEnergy
22             elif PhotonEnergy > secondPhotonEnergy:
23                 thirdPhotonIndex = secondPhotonIndex
24                 thirdPhotonEnergy = secondPhotonEnergy
25                 secondPhotonIndex = theParticle.index - 1
26                 secondPhotonEnergy = PhotonEnergy
27             elif PhotonEnergy > thirdPhotonEnergy:
28                 thirdPhotonIndex = theParticle.index - 1
29                 thirdPhotonEnergy = PhotonEnergy
30         if largestPhotonIndex >= 0 and secondPhotonIndex >= 0 and thirdPhotonIndex
           >= 0:
31             toAdd = theEvent.particles[largestPhotonIndex].momentum.values
32             toAdd = toAdd + theEvent.particles[secondPhotonIndex].momentum.values
33             toAdd = toAdd + theEvent.particles[thirdPhotonIndex].momentum.values
34             toAdd = toAdd + [tag]
35             result.append(toAdd)
36         idx = idx + 1
37     return result

```

The above code loads data from an LHCO file and converts the data to the format which is ready for the KMAD algorithm. The data set contains information about particle collisions involving the production of three photons.

The main function in the first code, “ChooseEventWithStrategy”, works by iterating over each event in “allEvents”, selecting the three photons with the largest energies, and combining their momenta into a single vector along with the tag. It keeps track of the largest, second-largest, and third-largest photons, and updates these values as it iterates over the particles in the event. If it finds an event with at least three photons, it adds the combined momentum vector to the result list. If result contains “count” elements, the function returns the result list.

```

1  import os
2
3  from Applications.kmeans.kmeansfunctions import ChooseEventWithStrategy,
    SaveCSVFile
4  from CutAndExport.CutEvent import CutEvents
5  from CutAndExport.CutFunctions import PhotonNumberCut
6  from Interfaces.LHCOlympics import LoadLHCOlympics
7
8  os.chdir("../..")
9
10
11 headList = ["FT0", "FT2", "FT5", "FT7", "FT8", "FT9"]
12 energyList = ["1500", "5000", "7000", "15000"]
13 PhotonNumberCut = PhotonNumberCut(1, [3])
14
15 for he in headList:
16     for en in energyList:

```

```

17     for i in range(0, 11):
18         testEvent = LoadLHCOlympics("_DataFolder/triphoton/cs/{0}/{0}-{1}-{2}.
            lhco".format(he, en, i))
19         CutEvents(testEvent, PhotonNumberCut)
20         resultList = ChooseEventWithStrategy(testEvent, len(testEvent.events),
            0)
21         toSave = "_DataFolder/kmeans/cs/E{0}/{1}/{1}-{0}-{2}.csv".format(en, he
            , i)
22         SaveCSVFile(toSave, resultList)
23         print(toSave, "saved! with events:", len(testEvent.events))

```

The above code imports various functions and libraries, including “os” for directory navigation and “LoadLHCOlympics” for loading the LHCO data set. It sets up some variables: “headList”, “energyList”, and “PhotonNumberCut”, in which “headList” and “energyList” are lists of strings representing different types of events in the LHCO files, while “PhotonNumberCut” is a Cut object used to filter events with fewer than 3 photons introduced in the previous section. The “PhotonNumberCut” is applied as the data cleaning phase. The strings in “headList” correspond to the origins of the NP signals, i.e. the  $O_{T_i}$  operators contributing to the aQGCs. The strings in “energyList” correspond to the beam energies of the muon colliders. The code loops over each combination of “headList”, “energyList”, and event index “i”. For each combination, it loads the corresponding LHCO file using “LoadLHCOlympics”, applies the “PhotonNumberCut” to filter out events with fewer than 3 photons using “CutEvents”, chooses events with the “ChooseEventWithStrategy” function, saves the results to a “csv” file using “SaveCSVFile”. Finally, the code prints the file name and the number of events saved.

After the “csv” files are saved, they can be directly used in the KMAD, IF and NIF algorithms to investigate the efficiency of the algorithms in searching for NP signals.

In the case of KMAD, `MLAnalysis` have functions “KMeans” and “CalculateDistance” in “Applications/kmeans/kmeansfunctions”,

```
def KMeans(dataList, d: int, k: int, nmin: int = 0) -> bool:
```

The “KMeans” function is fed with a data table as “dataList”, and number of features are provided as “d”, the number of clusters are provided as “k”. The data table is required to have one more column than the number of features which is used to store the cluster assignments of each point. There is an optional parameter “nmin”, which specifies a number of points (denoted as  $n_{min}$ ). When the number of points which change the cluster assignments is smaller than  $n_{min}$ , the function will stop. The function fails when there is at least a cluster without events. The returned value specifies whether the function succeeds. The cluster assignments are stored as the last elements of the vectors.

```
def CalculateDistance(dataList, d: int, k: int):
```

The input parameters of function “CalculateDistance” are the same as the “KMeans” function. The data table is required to have one more column than the number of features which stores the cluster assignments of the points. The returned value is an array storing the minimum distances for all points, where the minimum distance for a point is defined as the minimal distance between the point to all centroids.

The following code is an example to use the KMAD.

```

1  import numpy as np
2  from Applications.kmeans.kmeansfunctions import KMeans, CalculateDistance
3
4  k = 50

```

```

5  l = 100
6  dim = 12
7  data = np.loadtxt(data.csv", delimiter=',')
8  averageDistance=[]
9  for i in range(0, l):
10     succeed=False
11     while not succeed:
12         succeed=KMeans(data, dim, k)
13         np.savetxt("cluster-{0}.csv".format(i), data[:, dim].astype(int), delimiter
                    =',', fmt='%i')
14         distance=CalculateDistance(data, dim, k)
15         averageDistance.append(distance)
16     npAllDistance=np.array(averageDistance)
17     npAllDistance=np.transpose(npAllDistance)
18     npAverageDistance=np.mean(npAllDistance, axis=1)

```

The above code loads a data table from a “csv” file, and then use the “KMeans” function to apply the K-means algorithm. Then, the cluster assignments are stored for further usages. The minimum distances are calculated using “CalculateDistance”. Since the result depends on the randomly initialed cluster assignments, the code repeatedly calculates minimum distances for  $l = 100$  times and stores the results in “npAllDistance”. The average minimum distance “npAverageDistance” is then calculated as the anomaly scores.

The core of both the IF and NIF algorithm is the construction of an iTree. **MLAnalysis** has a function called “Split” in “Applications/IsolationForest/IsolationTree.py” to construct an iTree.

```
def Split(dataArray, length: int, maxSplit: int):
```

The “Split” function built an iTree with the data fed as “dataArray”, and number of features are provided as “length”. Two columns need to be added at the end of the data table to be used as results. The first of these columns is used to hold the source label of the data (e.g., whether it is from background or signal), and the function “Split” will not use this information or change the results in this column. The second column is used to hold the depth of the leaf in the iTree corresponding to a piece of data. The parameter “maxSplit” specifies the maximum depth of the leafs of an iTree. When “maxSplit” is set to  $-1$ , the iTree will be built with one point on one leaf. The returned values are the data sets with the last column set as the depths of the leafs.

The following code is an example to use the “Split” function to build iTrees.

```

1  Loop = 100
2  L = 18
3  saveCol = [L, L + 1]
4  dataSet1 = np.loadtxt("background.csv", delimiter=',')
5  dataSet2 = np.loadtxt("signal.csv", delimiter=',')
6  length = len(dataSet1)
7  z1 = np.zeros([length, 2])
8  length = len(dataSet2)
9  o2 = np.ones([length, 1])
10 z2 = np.zeros([length, 1])
11 dataSet1 = np.hstack((dataSet1, z1))
12 dataSet2 = np.hstack((dataSet2, o2, z2))
13 dataSet = np.vstack((dataSet1, dataSet2))
14 for i in range(0, Loop):
15     print("=====Loop{0}=====".format(i + 1))

```

```

16     resSet = Split(dataSet, L, -1)
17     np.savetxt(saveName + str(i) + ".csv", resSet[:, saveCol].astype(int),
                delimiter=',', fmt='%i')

```

The above code loads two data tables with  $L = 18$  features, and then use the “Split” function to build  $Loop = 100$  iTrees (i.e. an IF). Then, the tags and the depths of the leafs are stored in a “csv” file for further usages. For example, the average depths can be directly used as the anomaly scores in the IF algorithm.

#### 4.2. Using artificial neural network to reconstruct the energy of subprocess

For an effective field theory (EFT), the Wilson coefficients are usually functions of energy scales. Apart from that, the energy scale of the process is the basis for studying the validity of the EFT. However, the energy of the VBS subprocess can only be reconstructed by the final state particles. If there are multiple neutrinos in the final state, the traditional method based on kinematics is difficult to reconstruct the energy of the process. Artificial neural networks (ANN) are effective at finding complex relationships between inputs and outputs, which is suitable for solving such problems. Therefore, in the following example, we will present the data preparation phase of using the ANN to reconstruct the energy scale of a VBS subprocess with multiple neutrinos in the final state which is used in Ref. [21].

Before using ANN, we need training data sets and validation data sets. The LHE files created by MadGraph5\_aMC@NLO contain the information of neutrinos, so the  $\hat{s}$  can be calculated. With  $\hat{s}$  calculated, it is possible to build the label for the output layer of the ANN.

```

1  def SHatWWReal(eventSample: EventSample) -> float:
2      """
3      pW + pW = pL + pL + pnu + pnu
4      assume LHE file, so nu is visible
5      :param eventSample:
6      :return:
7      """
8      pAll = LorentzVector(0, 0, 0, 0)
9      for i in range(len(eventSample.particles)):
10         if 1 <= eventSample.particles[i].particleType <= 2:
11             pAll = pAll + eventSample.particles[i].momentum
12         if ParticleType.Missing == eventSample.particles[i].particleType:
13             pAll = pAll + eventSample.particles[i].momentum
14     return pAll * pAll

```

This code defines a function called “SHatWWReal” that takes an object of type “EventSample” as input and returns a float value. The function assumes that the input EventSample object is from an LHE file so that the neutrinos are visible.

The function calculates the square of the total momentum of the system resulting from the collision of two  $W$  bosons, which decay into two charged leptons and two neutrinos. The “pAll” variable is initialized as a Lorentz vector with zero momentum and iterates over all the particles in the given event sample. If the particle type is 1 or 2, which represent the charged leptons (Electron=1, Muon=2), then their momenta are added to the total four-momentum. If the particle type is “Missing”, which represents the neutrinos, then their momenta are also added to the total four-momentum. Finally, the function returns the square of the total momentum of the system.

The output of this function is used as a label for the output layer of an ANN that is trained to reconstruct the energy of a sub-process in a collision event. When ‘SHatWWReal’ function is implemented, we need to read an LHCO file and the corresponding LHE file to build the training and validation data sets.

```

1  def Export(eventSetLHCO, eventSetLHE, startIndex, endIndex, applyCut, file):
2      normalizer = 1000.0
3      # result_f.write("j1x,j1y,j1z,j2x,j2y,j2z,l1x,l1y,l1z,l2x,l2y,l2z,mx,my,shat\n"
4      #               ")
5      for i in range(startIndex, endIndex):
6          oneEvent = eventSetLHCO.events[i]
7          lepton1 = LorentzVector(0, 0, 0, 0)
8          lepton2 = LorentzVector(0, 0, 0, 0)
9          jet1 = LorentzVector(0, 0, 0, 0)
10         jet2 = LorentzVector(0, 0, 0, 0)
11         missing = LorentzVector(0, 0, 0, 0)
12         largestJetIndex1 = 0
13         largestJetM1 = 0.0
14         largestJetIndex2 = 0
15         largestJetM2 = 0.0
16         leptonIdx1 = 0
17         leptonIdx2 = 0
18         largestLepton1 = 0
19         largestLepton2 = 0
20         hasMissing = False
21         for oneParticle in oneEvent.particles:
22             if ParticleStatus.Outgoing == oneParticle.status \
23                 and ParticleType.Jet == oneParticle.particleType:
24                 momentum = oneParticle.momentum.Momentum()
25                 if momentum > largestJetM1:
26                     largestJetM2 = largestJetM1
27                     largestJetIndex2 = largestJetIndex1
28                     largestJetM1 = momentum
29                     largestJetIndex1 = oneParticle.index
30                 elif momentum > largestJetM2:
31                     largestJetM2 = momentum
32                     largestJetIndex2 = oneParticle.index
33             elif ParticleType.Electron <= oneParticle.particleType <= ParticleType.
34                 Muon:
35                 momentumLepton = oneParticle.momentum.Momentum()
36                 if oneParticle.PGDid > 0 and momentumLepton > largestLepton1:
37                     leptonIdx1 = oneParticle.index
38                     largestLepton1 = momentumLepton
39                 elif oneParticle.PGDid < 0 and momentumLepton > largestLepton2:
40                     leptonIdx2 = oneParticle.index
41                     largestLepton2 = momentumLepton
42             elif ParticleType.Missing == oneParticle.particleType:
43                 hasMissing = True
44                 missing = missing + oneParticle.momentum
45         if not (leptonIdx1 > 0 and leptonIdx2 > 0):
46             continue
47         if not (largestJetIndex1 > 0 and largestJetIndex2 > 0):
48             continue
49         if not hasMissing:
50             print(oneEvent.DebugPrint())
51             continue
52         lepton1 = oneEvent.particles[leptonIdx1 - 1].momentum
53         lepton2 = oneEvent.particles[leptonIdx2 - 1].momentum

```

```

52     jet1 = oneEvent.particles[largestJetIndex1 - 1].momentum
53     jet2 = oneEvent.particles[largestJetIndex2 - 1].momentum
54     if applyCut:
55         if lepton1.values[1] * lepton1.values[1] + lepton1.values[2] * lepton1.
           values[2] < 1.0e2:
56             continue
57         if lepton2.values[1] * lepton2.values[1] + lepton2.values[2] * lepton2.
           values[2] < 1.0e2:
58             continue
59         if missing.values[1] * missing.values[1] + missing.values[2] * missing.
           values[2] < 1.0e2:
60             continue
61         lepX = lepton1.values[1] + lepton2.values[1]
62         lepY = lepton1.values[2] + lepton2.values[2]
63         if lepX * lepX + lepY * lepY < 100:
64             continue
65         lengthLep = math.sqrt(lepX * lepX + lepY * lepY)
66         lengthM = math.sqrt(missing.values[1] * missing.values[1] + missing.
           values[2] * missing.values[2])
67         dotLM = (lepX * missing.values[1] + lepY * missing.values[2]) / (
           lengthLep * lengthM)
68         if abs(dotLM) < 0.8:
69             continue
70         lengthL1 = math.sqrt(lepton1.values[1] * lepton1.values[1]
           + lepton1.values[2] * lepton1.values[2]
           + lepton1.values[3] * lepton1.values[3])
71         lengthL2 = math.sqrt(lepton2.values[1] * lepton2.values[1]
           + lepton2.values[2] * lepton2.values[2]
           + lepton2.values[3] * lepton2.values[3])
72         dotLL = (lepton1.values[1] * lepton2.values[1]
           + lepton1.values[2] * lepton2.values[2]
           + lepton1.values[3] * lepton2.values[3]) / (lengthL1 * lengthL2
           )
73         if dotLL > -0.5:
74             continue
75     realShat = SHatWWReal(eventSetLHE.events[i])
76     paramLst = [jet1.values[0] / normalizer,
77                 jet1.values[1] / normalizer,
78                 jet1.values[2] / normalizer,
79                 jet1.values[3] / normalizer,
80                 jet2.values[0] / normalizer,
81                 jet2.values[1] / normalizer,
82                 jet2.values[2] / normalizer,
83                 jet2.values[3] / normalizer,
84                 lepton1.values[0] / normalizer,
85                 lepton1.values[1] / normalizer,
86                 lepton1.values[2] / normalizer,
87                 lepton1.values[3] / normalizer,
88                 lepton2.values[0] / normalizer,
89                 lepton2.values[1] / normalizer,
90                 lepton2.values[2] / normalizer,
91                 lepton2.values[3] / normalizer,
92                 missing.values[1] / normalizer,
93                 missing.values[2] / normalizer]
94     strW = ""
95     for x in range(0, 18):
96         strW = "{}{:.5e},".format(strW, paramLst[x])
97     strW = "{}{:.5e}\n".format(strW, math.sqrt(realShat) / normalizer)
98     file.write(strW)

```

---

The code above defines a function called “Export” that takes in six parameters: “eventSetLHCO”, “eventSetLHE”, “startIndex”, “endIndex”, “applyCut”, and “file”. Inside the function, a variable normalizer is initialized with a value of 1000.0, the goal is to change the unit from GeV to TeV.

The data in the above code is arranged as a  $19 \times N$  table, where  $N$  is the number of events, each line is a 19 dimensional vector represents an event. They are the components of the 4-momenta of the two hardest jets, the 4-momenta of the two hardest opposite signed charged leptons and the components of the transverse missing momentum, they are all observables. The output of the ANN corresponds to  $\hat{s}$ . The true labels are the 19-th variables of the elements in the data sets which are  $\hat{s}_{\text{tr}}$  of the events.

At the beginning of the code, we iterate through each particle in the event, selecting the two jets with the highest momenta and the two leptons (either electrons or muons) with the highest momenta. We require that these two leptons are a pair of oppositely charged leptons. Additionally, the function checks if there is any missing momentum in the event. It can be found that not all events are useful, such as the events where two jets are not found, events where two leptons with the largest momenta are not a pair of opposite-charged particles, and events where missing energy is not found. Therefore, to ensure the effectiveness of the subsequent code, we perform data cleaning on events.

In this study, the NP is introduced via the SMEFT. As an EFT, the  $\hat{s}$  is important. However, we do not need to reconstruct the  $\hat{s}$  of the SM events. To increase the efficiency of the ANN, we can remove the SM events. We introduce a switch called “applyCut” in the second part of the code, which includes a series of filtering strategies: requiring the transverse momentum (“pT”) of each charged lepton to be greater than 100 GeV, requiring the missing transverse momentum to be greater than 100 GeV, requiring the transverse momentum of the lepton pair combination to be greater than 100 GeV, requiring the cosine of the angle between the transverse momentum of the lepton pair and the missing transverse momentum to be greater than 0.8, and requiring the cosine of the angle between the pair of oppositely charged leptons to be less than -0.5. If an event passes all the selection criteria, it is more likely from the NP and is written to the output file. Those cuts are the cuts to high light the NP contribution used in Ref. [27]. When “applyCut” is set to true, the occurrence of background events can be suppressed, resulting in training and validation data sets that mainly contain signal events.

The end part of the code writes the selected events and their corresponding input parameters to a file. Afterwards, the parameters are appended to the string “strW” using Python’s “format” function and written to the file using Python’s “write” function. Each line of the file corresponds to an event, with the input parameters and output values separated by commas.

Finally, the data is saved as a “csv” file, and is ready to be fed to the ANN [21, 38].

#### 4.3. ESS using polarization

In order to distinguish between the SM events and NP events, one can utilize certain special properties such as polarization. Take the VBS process as an example, for many events, the charged leptons in the final states are from the  $Z$  boson decays, whose polarizations can be transverse or longitudinal. In the rest frame of a  $Z$  boson with the flight

direction of the  $Z$  boson as the  $\mathbf{z}$ -axis, the charged leptons in the  $Z$  boson decay satisfy an angular distribution [39],

$$\frac{d\sigma}{d\cos\theta^*} \propto f_L \frac{(1 - \cos\theta^*)^2}{4} + f_R \frac{(1 + \cos\theta^*)^2}{4} + f_0 \frac{\sin^2\theta^*}{2}, \quad (1)$$

where  $f_{L,R,0}$  represent the fraction of left, right-handed and longitudinal polarizations, and  $\theta^*$  is the zenith angle of the charged lepton in the rest frame of  $Z$  boson.

In the SM,  $Z$  bosons are mainly longitudinally polarized. In some NP models,  $Z$  bosons are mainly transversely polarized, leading to a different angular distribution. Therefore, we can extract NP signal events from the SM background by utilizing the angular distribution of the charged leptons produced by the decay of the  $Z$  boson in its rest frame.

In the code below, we define a function called “leptonPZ”. The first 23 lines of the function traverse all particles in the event and select the two leptons with the largest momenta (“largestLM1” and “largestLM2”). Then, the sum of the momenta of the two selected leptons is assigned to “pZ”, and the momentum of the positive charged lepton is assigned to “pL”. That is, we reconstruct the momentum of the  $Z$  boson (“pZ”) from the sum of the momenta of the lepton pair produced by its decay.

In the latter part of the code, we use two functions in “Matrix4x4”. First, we use the “MakeRotation” function to create a rotation matrix (“rotMtr”) that rotates the momentum of the  $Z$  boson (“pZ”) such that it aligns with the  $\mathbf{z}$ -axis. It then applies the same rotation to the momentum of the lepton (“pL”). This rotation matrix acts as a rotation of the frame, and “pLDir” is the result of the momentum of the charged lepton after the rotation is applied, where the rotation is the one rotates the momentum of the  $Z$  boson to the  $\mathbf{z}$ -axis. Then, we use the “MakeBoost” function to generate a Lorentz transformation matrix that transforms “pLDir” into the rest frame of the  $Z$  boson, obtaining “pLRest”. Finally, the function returns the cosine of the zenith angle of “pLRest” in the rest frame of the  $Z$  boson (the  $\cos(\theta^*)$  in Eq. (1)). In this way, one can apply a cut on  $\cos(\theta^*)$  to highlight the NP signals if the polarizations of the gauge bosons are different from the SM [28–30].

```

1  def LeptonPZ(eventSample: EventSample) -> float:
2      largestLM1 = 0
3      largestLM2 = 0
4      largestLIndex1 = 0
5      largestLIndex2 = 0
6      for particle in eventSample.particles:
7          if ParticleType.Electron == particle.particleType or ParticleType.Muon ==
            particle.particleType:
8              momentum = particle.momentum.Momentum()
9              if momentum > largestLM1:
10                  largestLM2 = largestLM1
11                  largestLIndex2 = largestLIndex1
12                  largestLM1 = momentum
13                  largestLIndex1 = particle.index
14              elif momentum > largestLM2:
15                  largestLM2 = momentum
16                  largestLIndex2 = particle.index
17      p41 = eventSample.particles[largestLIndex1 - 1].momentum
18      p42 = eventSample.particles[largestLIndex2 - 1].momentum
19      pZ = p41 + p42

```

```

20     pL = p41
21     if eventSample.particles[largestLIndex2 - 1].PGDid > 0:
22         pL = p42
23         rotMtr = Matrix4x4.MakeRotationFromTo(pZ.V3d(), [0, 0, 1])
24         pZDir = rotMtr.MultiplyVector(pZ)
25         pLDir = rotMtr.MultiplyVector(pL)
26         v3dz = pZDir.V3d()
27         vsq = Constants.dot3d(v3dz, v3dz)
28         if vsq > 0.9999999999 or vsq < 0:
29             # This event sample should be excluded
30             return 1000
31         boostMtr = Matrix4x4.MakeBoost(pZDir.V3d())
32         pLRest = boostMtr.MultiplyVector(pLDir)
33         return math.cos(pLRest.Theta())

```

## 5. Summary

ML algorithms are widely used in many fields of HEP. **MLAnalysis** provides a suite of tools for data transformation and feature engineering. It converts experimental or simulated data into a data set that can be used for ML, through which a variety of algorithms such as ANN can be easily applied. The current version of the **MLAnalysis** has built-in ML algorithms including IF, NIF and KMAD. These ML-based approaches can help researchers to optimize ESS to improve signal significance.

In Sec. 3, we detail the IF, NIF and KMAD algorithms and their applications. IF and KMAD algorithms are AD algorithms which are NP model-independent and SMEFT operator-independent. No matter what the NP signals exist in the data set, they can be detected as long as they satisfy “few and different” from the SM events. NIF algorithm not only inherits the advantages of IF algorithm, but also can deal with non-AD problems, such as detecting NP effects dominated by interference. As automatic ESSs, IF, NIF and KMAD can achieve signal recognition ability comparable to or even better than that of a traditional ESS without kinematic analysis, and greatly improve the analysis efficiency without knowing what kind of NP model or SMEFT operator one is studying. We presented several examples of **MLAnalysis** in Sec. 4, which include how to build a data set for ML, how to reconstruct the energy scale of the subprocess and how to optimize the ESS using polarization.

The combination of data preparation and ESSs in **MLAnalysis** makes it a comprehensive and valuable tool for researchers working with experimental and MC simulation data. It has the potential to make ML more accessible and effective in the NP research. If users want to apply ML algorithms to other problems beyond the scope of this paper, the open-source **MLAnalysis** based on user-friendly Python code enables a path forward for exploration.

## Conflict of interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgment

This work was supported in part by the National Natural Science Foundation of China under Grants Nos. 11905093 and 12147214, the Natural Science Foundation of the Liaoning Scientific Committee No. LJKZ0978 and the Outstanding Research Cultivation Program of Liaoning Normal University (No.21GDL004) and “New strategies for detecting signal of new physics at the future lepton colliders”, and the University-Industry Collaborative Education Program No. 220800575313412.

## References

- [1] A. Radovic, M. Williams, D. Rousseau, M. Kagan, D. Bonacorsi, A. Himmel, A. Aurisano, K. Terao, T. Wongjirad, Machine learning at the energy and intensity frontiers of particle physics, *Nature* 560 (7716) (2018) 41–48. [doi:10.1038/s41586-018-0361-2](#).
- [2] P. Baldi, P. Sadowski, D. Whiteson, Searching for Exotic Particles in High-Energy Physics with Deep Learning, *Nature Commun.* 5 (2014) 4308. [arXiv:1402.4735](#), [doi:10.1038/ncomms5308](#).
- [3] J. Ren, L. Wu, J. M. Yang, J. Zhao, Exploring supersymmetry with machine learning, *Nucl. Phys. B* 943 (2019) 114613. [arXiv:1708.06615](#), [doi:10.1016/j.nuclphysb.2019.114613](#).
- [4] M. Abdughani, J. Ren, L. Wu, J. M. Yang, Probing stop pair production at the LHC with graph neural networks, *JHEP* 08 (2019) 055. [arXiv:1807.09088](#), [doi:10.1007/JHEP08\(2019\)055](#).
- [5] A. De Simone, T. Jacques, Guiding New Physics Searches with Unsupervised Learning, *Eur. Phys. J. C* 79 (4) (2019) 289. [arXiv:1807.06038](#), [doi:10.1140/epjc/s10052-019-6787-3](#).
- [6] J. Ren, L. Wu, J. M. Yang, Unveiling CP property of top-Higgs coupling with graph neural networks at the LHC, *Phys. Lett. B* 802 (2020) 135198. [arXiv:1901.05627](#), [doi:10.1016/j.physletb.2020.135198](#).
- [7] R. T. D’Agnolo, A. Wulzer, Learning New Physics from a Machine, *Phys. Rev. D* 99 (1) (2019) 015014. [arXiv:1806.02350](#), [doi:10.1103/PhysRevD.99.015014](#).
- [8] M. van Beekveld, S. Caron, L. Hendriks, P. Jackson, A. Leinweber, S. Otten, R. Patrick, R. Ruiz De Austri, M. Santoni, M. White, Combining outlier analysis algorithms to identify new physics at the LHC, *JHEP* 09 (2021) 024. [arXiv:2010.07940](#), [doi:10.1007/JHEP09\(2021\)024](#).
- [9] R. Iten, T. Metger, H. Wilming, L. del Rio, R. Renner, [Discovering physical concepts with neural networks](#), *Phys. Rev. Lett.* 124 (2020) 010508. [arXiv:1807.10300](#), [doi:10.1103/PhysRevLett.124.010508](#).  
URL <https://link.aps.org/doi/10.1103/PhysRevLett.124.010508>
- [10] M. Crispim Romão, N. F. Castro, R. Pedro, Finding New Physics without learning about it: Anomaly Detection as a tool for Searches at Colliders, *Eur. Phys. J. C* 81 (1) (2021) 27, [Erratum: *Eur.Phys.J.C* 81, 1020 (2021)]. [arXiv:2006.05432](#), [doi:10.1140/epjc/s10052-021-09813-2](#).
- [11] E. Fol, R. Tomás, J. Coello de Portugal, G. Franchetti, Detection of faulty beam position monitors using unsupervised learning, *Phys. Rev. Accel. Beams* 23 (10) (2020) 102805. [doi:10.1103/PhysRevAccelBeams.23.102805](#).
- [12] M. A. Md Ali, N. Badrud’din, H. Abdullah, F. Kemi, Alternate methods for anomaly detection in high-energy physics via semi-supervised learning, *Int. J. Mod. Phys. A* 35 (23) (2020) 2050131. [doi:10.1142/S0217751X20501316](#).
- [13] H. Lv, D. Wang, L. Wu, Deep learning jet images as a probe of light Higgsino dark matter at the LHC, *Phys. Rev. D* 106 (5) (2022) 055008. [arXiv:2203.14569](#), [doi:10.1103/PhysRevD.106.055008](#).
- [14] A. J. Larkoski, I. Moult, B. Nachman, Jet Substructure at the Large Hadron Collider: A Review of Recent Advances in Theory and Machine Learning, *Phys. Rept.* 841 (2020) 1–63. [arXiv:1709.04464](#), [doi:10.1016/j.physrep.2019.11.001](#).
- [15] J. Guo, J. Li, T. Li, F. Xu, W. Zhang, Deep learning for  $R$ -parity violating supersymmetry searches at the LHC, *Phys. Rev. D* 98 (7) (2018) 076017. [arXiv:1805.10730](#), [doi:10.1103/PhysRevD.98.076017](#).
- [16] M. Abdughani, J. Ren, L. Wu, J. M. Yang, J. Zhao, Supervised deep learning in high energy phenomenology: a mini review, *Commun. Theor. Phys.* 71 (8) (2019) 955. [arXiv:1905.06047](#), [doi:10.1088/0253-6102/71/8/955](#).
- [17] J. Li, S. Yang, R. Zhang, Detecting anomalies in vector boson scattering, *Chin. Phys. C* 45 (7) (2021) 073104. [arXiv:2010.13281](#), [doi:10.1088/1674-1137/abf829](#).

- [18] G. Kasieczka, et al., The LHC Olympics 2020 a community challenge for anomaly detection in high energy physics, Rept. Prog. Phys. 84 (12) (2021) 124201. [arXiv:2101.08320](#), [doi:10.1088/1361-6633/ac36b9](#).
- [19] Y.-C. Guo, L. Jiang, J.-C. Yang, Detecting anomalous quartic gauge couplings using the isolation forest machine learning algorithm, Phys. Rev. D 104 (3) (2021) 035021. [arXiv:2103.03151](#), [doi:10.1103/PhysRevD.104.035021](#).
- [20] J.-C. Yang, Y.-C. Guo, L.-H. Cai, Using a nested anomaly detection machine learning algorithm to study the neutral triple gauge couplings at an  $e^+e^-$  collider, Nucl. Phys. B 977 (2022) 115735. [arXiv:2111.10543](#), [doi:10.1016/j.nuclphysb.2022.115735](#).
- [21] J.-C. Yang, J.-H. Chen, Y.-C. Guo, Extract the energy scale of anomalous  $\gamma\gamma \rightarrow W^+W^-$  scattering in the vector boson scattering process using artificial neural networks, JHEP 09 (2021) 085. [arXiv:2107.13624](#), [doi:10.1007/JHEP09\(2021\)085](#).
- [22] J.-C. Yang, X.-Y. Han, Z.-B. Qin, T. Li, Y.-C. Guo, Measuring the anomalous quartic gauge couplings in the  $W^+W^- \rightarrow W^+W^-$  process at muon collider using artificial neural networks, JHEP 09 (2022) 074. [arXiv:2204.10034](#), [doi:10.1007/JHEP09\(2022\)074](#).
- [23] J. Searcy, L. Huang, M.-A. Pleier, J. Zhu, Determination of the  $WW$  polarization fractions in  $pp \rightarrow W^\pm W^\pm jj$  using a deep machine learning technique, Phys. Rev. D 93 (9) (2016) 094033. [arXiv:1510.01691](#), [doi:10.1103/PhysRevD.93.094033](#).
- [24] J. Lee, N. Chanon, A. Levin, J. Li, M. Lu, Q. Li, Y. Mao, Polarization fraction measurement in same-sign  $WW$  scattering using deep learning, Phys. Rev. D 99 (3) (2019) 033004. [arXiv:1812.07591](#), [doi:10.1103/PhysRevD.99.033004](#).
- [25] J. Lee, N. Chanon, A. Levin, J. Li, M. Lu, Q. Li, Y. Mao, Polarization fraction measurement in  $ZZ$  scattering using deep learning, Phys. Rev. D 100 (11) (2019) 116010. [arXiv:1908.05196](#), [doi:10.1103/PhysRevD.100.116010](#).
- [26] J. Li, C. Zhang, R. Zhang, Polarization measurement for the dileptonic channel of  $W+W^-$  scattering using generative adversarial network, Phys. Rev. D 105 (1) (2022) 016005. [arXiv:2109.09924](#), [doi:10.1103/PhysRevD.105.016005](#).
- [27] Y.-C. Guo, Y.-Y. Wang, J.-C. Yang, Constraints on anomalous quartic gauge couplings by  $\gamma\gamma \rightarrow W^+W^-$  scattering, Nucl. Phys. B 961 (2020) 115222. [arXiv:1912.10686](#), [doi:10.1016/j.nuclphysb.2020.115222](#).
- [28] Y.-C. Guo, Y.-Y. Wang, J.-C. Yang, C.-X. Yue, Constraints on anomalous quartic gauge couplings via  $W\gamma jj$  production at the LHC, Chin. Phys. C 44 (12) (2020) 123105. [arXiv:2002.03326](#), [doi:10.1088/1674-1137/abb4d2](#).
- [29] Q. Fu, J.-C. Yang, C.-X. Yue, Y.-C. Guo, The study of neutral triple gauge couplings in the process  $e^+e^- \rightarrow Z\gamma$  including unitarity bounds, Nucl. Phys. B 972 (2021) 115543. [arXiv:2102.03623](#), [doi:10.1016/j.nuclphysb.2021.115543](#).
- [30] J.-C. Yang, Y.-C. Guo, C.-X. Yue, Q. Fu, Constraints on anomalous quartic gauge couplings via  $Z\gamma jj$  production at the LHC, Phys. Rev. D 104 (3) (2021) 035015. [arXiv:2107.01123](#), [doi:10.1103/PhysRevD.104.035015](#).
- [31] J.-C. Yang, Z.-B. Qing, X.-Y. Han, Y.-C. Guo, T. Li, Tri-photon at muon collider: a new process to probe the anomalous quartic gauge couplings, JHEP 22 (2022) 053. [arXiv:2204.08195](#), [doi:10.1007/JHEP07\(2022\)053](#).
- [32] F. T. Liu, K. M. Ting, Z. Zhou, Isolation forest, in: 2008 Eighth IEEE International Conference on Data Mining, 2008, pp. 413–422. [doi:10.1109/ICDM.2008.17](#).
- [33] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, O. Mattelaer, H. S. Shao, T. Stelzer, P. Torrielli, M. Zaro, The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations, JHEP 07 (2014) 079. [arXiv:1405.0301](#), [doi:10.1007/JHEP07\(2014\)079](#).
- [34] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna, S. Prestel, C. O. Rasmussen, P. Z. Skands, An introduction to PYTHIA 8.2, Comput. Phys. Commun. 191 (2015) 159–177. [arXiv:1410.3012](#), [doi:10.1016/j.cpc.2015.01.024](#).
- [35] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, M. Selvaggi, DELPHES 3, A modular framework for fast simulation of a generic collider experiment, JHEP 02 (2014) 057. [arXiv:1307.6346](#), [doi:10.1007/JHEP02\(2014\)057](#).
- [36] S. P. Lloyd, Least squares quantization in pcm, IEEE Trans. Inf. Theory 28 (1982) 129–136.
- [37] S. Zhang, J.-C. Yang, Y.-C. Guo, Using k-means assistant event selection strategy to study anomalous quartic gauge couplings at muon colliders (2 2023). [arXiv:2302.01274](#).
- [38] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems (3

- 2016). [arXiv:1603.04467](#).
- [39] M. Peruzzi, First measurement of vector boson polarization at LHC, Ph.D. thesis, Zurich, ETH (2011).