

HEP ML Lab

An end-to-end framework for machine learning application in high energy physics

Jing Li, Hao Sun

Dalian University of Technology

28th Mini-workshop on the frontier of LHC

Contents

Introduction

Pain points and motivation.

Core APIs

Most useful programming cases.

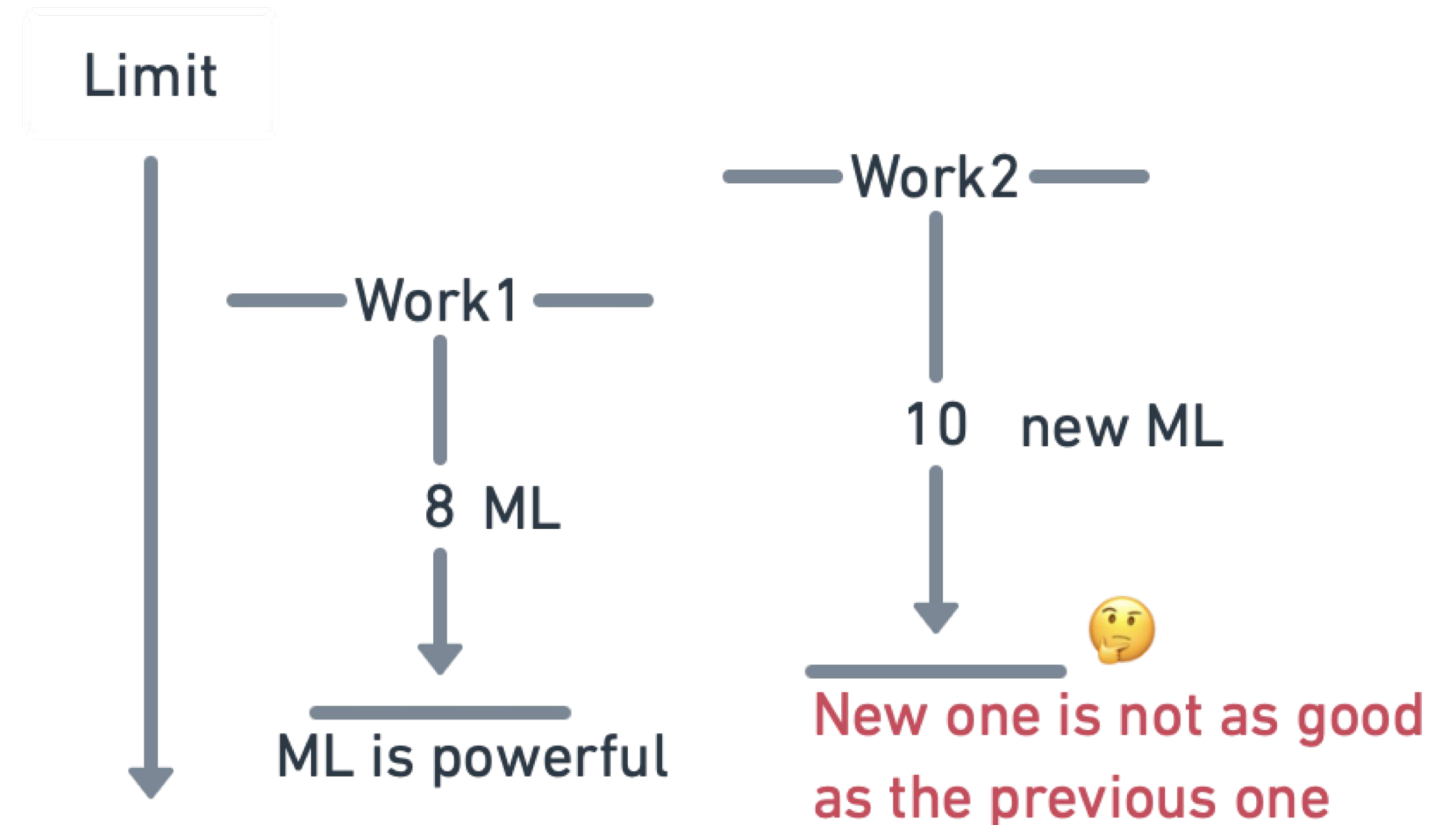
Roadmap

What's the progress and the next.

Introduction

Reproduction Issues

- A large amount of work explores the performance improvement brought by machine learning methods. Results are **promising** for the new physics search.
- However the lack of source codes makes it **quite difficult to reproduce** the results.
- If they are generated under different conditions, at what extent we can say that new methods are truly powerful and worth to try in broader subjects?



Introduction

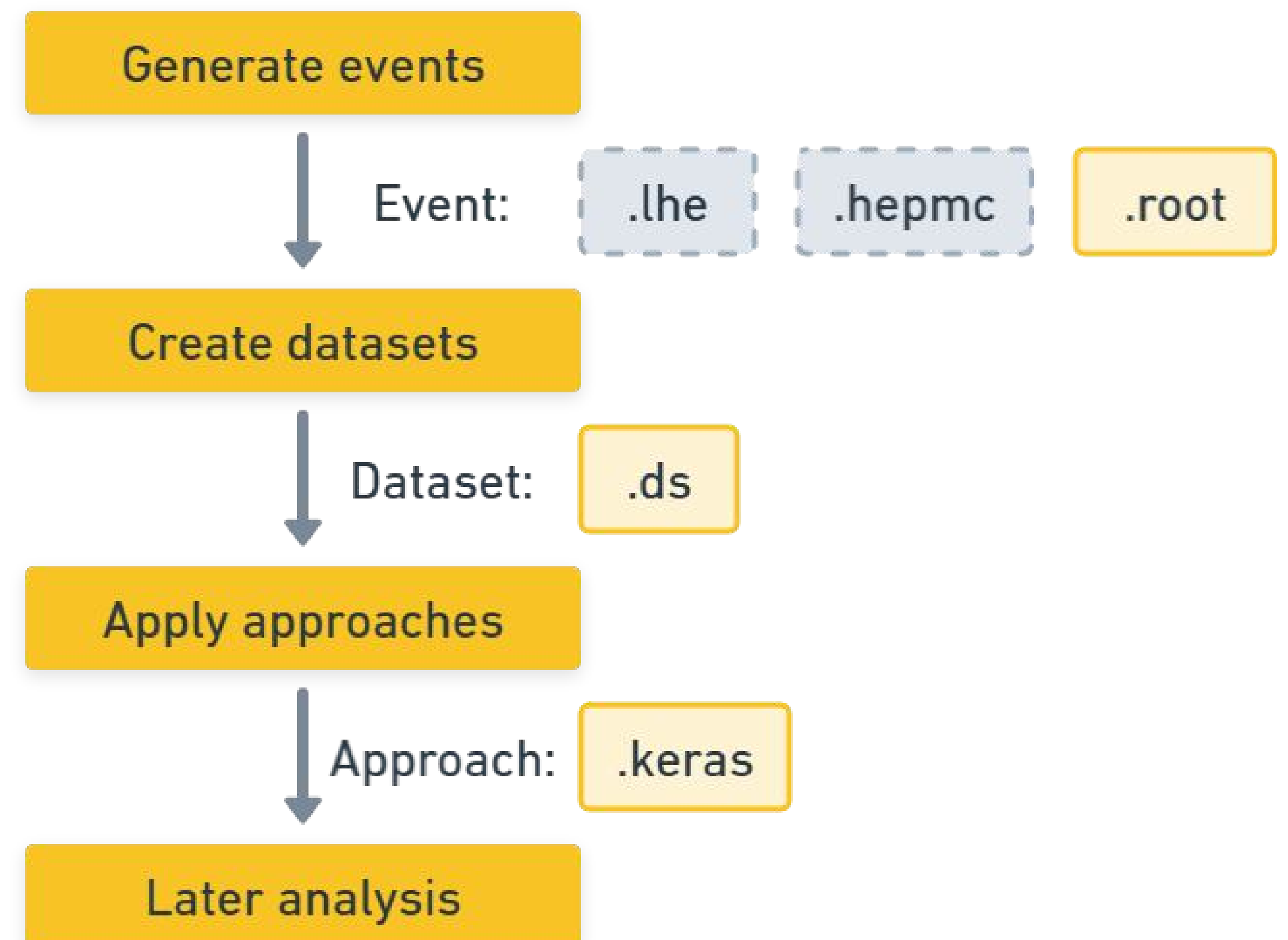
Existing frameworks

Name	Data	Model	Style	Highlight
<u>hep_ml</u>	No	Yes	sklearn	Low correlation, theano-based, sklearn compatible
<u>weaver</u>	Yes	Yes	CLI + config	Support many dataset formats, config for all
<u>JetNet</u>	Yes	No	custom	Three datasets, generative evaluation metrics
<u>pd4ml</u>	Yes	Yes	modified Keras	Five datasets, model templates
<u>MLAnalysis</u>	Yes	Yes	custom	LHE/LHCO data, three ML algorithms
<u>mapyde</u>	Yes	No	CLI + TUI + config	Madgraph5 workflow, for a specific problem
<u>madminer</u>	Yes	Yes	custom	Madgraph5 workflow, for a specific problem
<u>hep-ml-lab</u>	Yes	Yes	Keras	Madgraph5 workflow, not specific for one problem

Introduction

Overall workflow

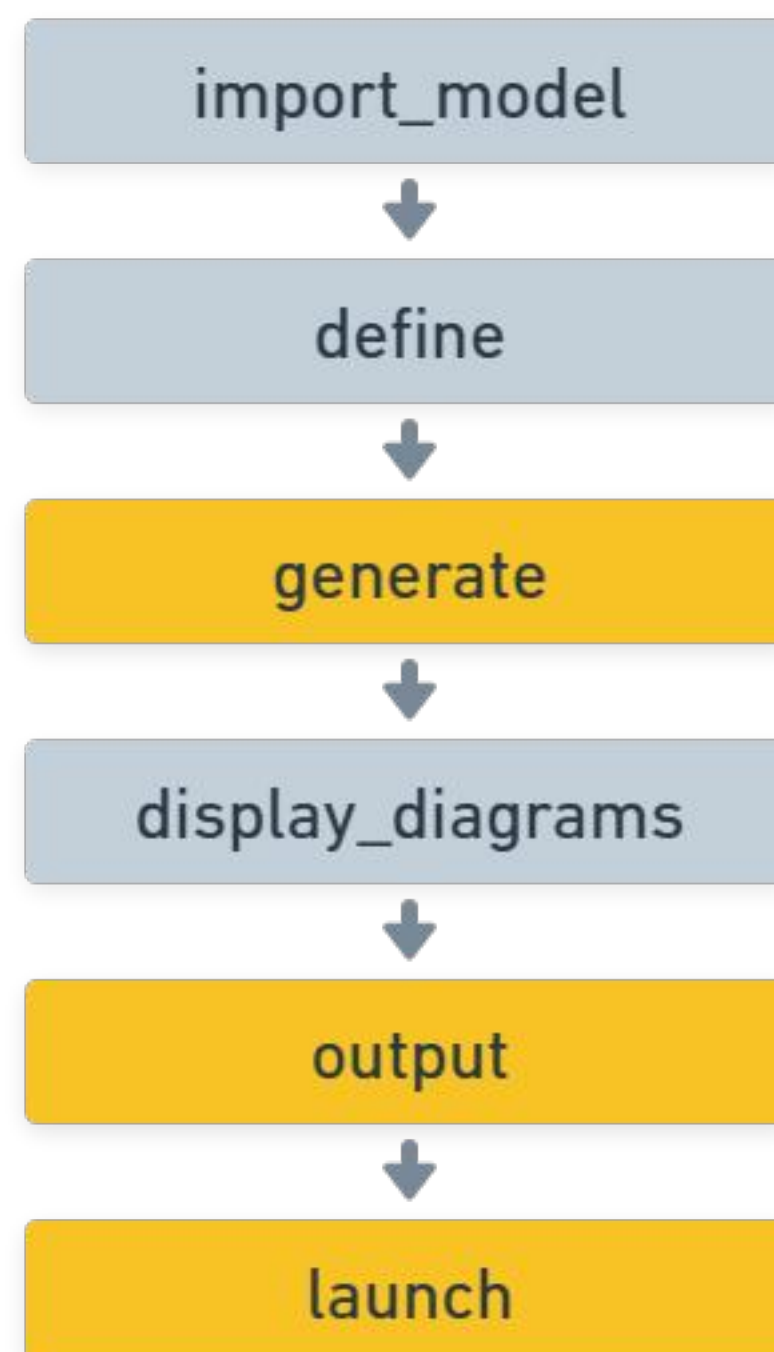
- An **end-to-end** framework for applying machine learning into HEP studies.
- **Simplify the data flow**: easier to manage and track different data.
- **Unify programming style** across all stages: objected-oriented and **Keras**.
- All makes the results more reliable and reproducible.



APIs

Generate events

- Minimal wrapper of Madgraph5



Made with  Whimsical

```
1 g = Madgraph5(executable="mg5_aMC", verbose=1)
2 g.generate(["p p > w+ z"])
3 g.output("data/pp2wz")
4 g.launch(
5     shower="pythia8",
6     detector="delphes",
7     madspin="none",
8     settings={
9         "nevents": 1000,
10        "run_tag": "250-300",
11        "pt_min_pdg": {24: 250},
12        "pt_max_pdg": {24: 300},
13    },
14    decays=[
15        "w+ > j j",
16        "z > vl vl~",
17    ],
18    seed=42,
19 )
```

APIs

Generate events

- Flexible control of intermediate output.
 - verbose=0: quiet
 - verbose=1: suppressed output as shown
 - verbose=2: original output
- Check the results via **summary**

```
Madgraph5_aMC@NLO v3.4.2
Total: 2 processes with 6 diagrams
1 diagrams saved in ./data/pp2wz/Diagrams
Output saved in ./data/pp2wz
run_01...survey...pythia8...delphes...storing...✓
```

p p > w+ z

#	Name	Collider	Tag	Cross section (pb)	N events	Seed
0	run_01[1]	pp:6500.0x6500.0	250-300	4.371e-02 +- 4.200e-04	1,000	42

Output: data/pp2wz

APIs

Events

Parton-level

.lhe

```
7000e-03 4.57238200e+02 7.54677100e-03 1.02559300e-03  
} 0 501 0 +0.0000000000e+00 +0.0000000000e+00  
} 0 0 501 -0.0000000000e+00 -0.0000000000e+00  
L 2 0 0 +4.2998352586e+02 +1.2575170009e+01  
L 2 0 0 -4.2998352586e+02 -1.2575170009e+01  
} 3 502 0 +1.6794615283e+00 -7.7917042810e-01  
} 3 0 502 +4.2830406433e+02 +1.3354340437e+01  
} 4 0 0 -4.2961266050e+02 -1.1631754951e+01  
L 4 0 0 -3.7086536509e-01 -9.4341505808e-02
```

- particle.pid
- particle.eta
- particle.phi
- particle.mass

Particle-level

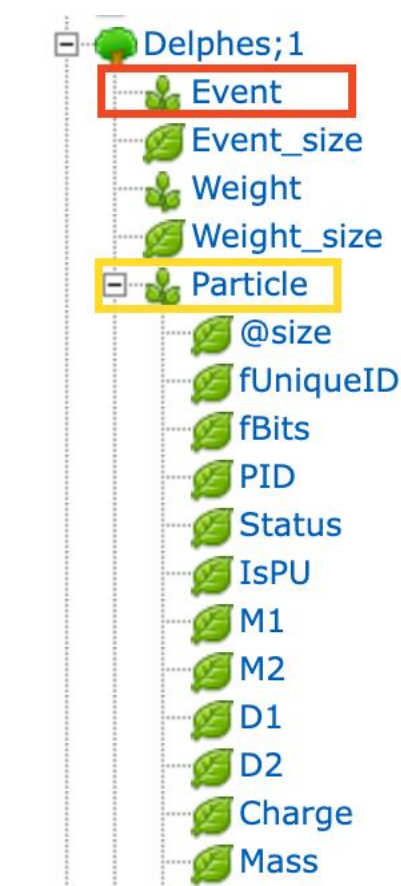
.hepmc

```
0000000000000000e+00 -1.0000000000000000e+00 -1.0000000000000000e+00  
1" "AUX_10" "AUX_100" "AUX_101" "AUX_102" "AUX_103" "AUX_104"  
999999998e-05 2.28370000000000001e-03  
0 0 0 1 0  
.6266581997519438e+02 6.6266581997519438e+02 0 21 0 0 -3 1  
0 0 0 2 0  
-4.0087158015244444e+02 4.0087158015244444e+02 0 21 0 0 -3  
3916869591523628e+01 1.2575872785131963e+01 -1.39651736562  
0 0 0 2 0
```

- particle.pid
- particle.eta
- particle.phi
- particle.mass

Detector-level

.root



- Electron.PT
- Jet.PT
- Muon.Charge
- Tower.ET
- Track.Mass
- ...

APIs

Events

```
1 events = R00TEvents.load("../..../tag_1_delphes_events.root")
2 events["jet.pt"]
```

```
[[520, 42.4],
```

```
...,
```

```
[545]]
```

```
-----  
type: 100 * var * float32
```

- Data is recorded as an awkward array.
- It could be interpreted as 100 events with variable number of jets.

APIs

Naming convention: PhysicsObject

- Physics object name = <type/alias>[<index>]

All muons.

muon

The **first 200** constituents of the **leading** jet.

jet0.constituents:200

The **first** subjet clustered by all fatjets' constituents with **kt** algorithm and **0.3** radius.

kt3fatjet0

APIs

Naming convention: Observable

- Observable name = <physics object>.<type/alias>

The charge of **all** muons's .

muon.charge

The transverse momentum of the **first 200** constituents of the **leading** jet.

jet0.constituents:200.pt

The distance between the **first** and the **second** subjects defined as before.

kt3jet0,kt3jet1.delta_r

→ The number of physics objects is determined by the observable.

APIs

Naming convention

```
1 events = load_events("../..../tag_1_delphes_events.root")
2
3 my_object = parse_physics_object("ak8jet.constituents")
4 my_object.read(events)
5 my_object.array
```

- Agnostic usage to ignore the event type.
- The familiar string, not programming classes.

```
1 parse_physics_object("ak8jet.constituents").read(events).array
```

- Method chaining makes it more compact.
- recorded observables

⇒ **100** * **var** * **var** * {
events jets constituents **pt: float32,**
number number number **eta: float32,**
phi: float32,
mass: float32
}

APIs

Naming convention

- Users can **register** aliases for physics objects.

```
1  subjects = parse_physics_object("ak3jet")
2  register(subjects, "subjects")
3
4  # In another place, get back the subjects quickly
5  parsed = parse_physics_object("subjects")
```

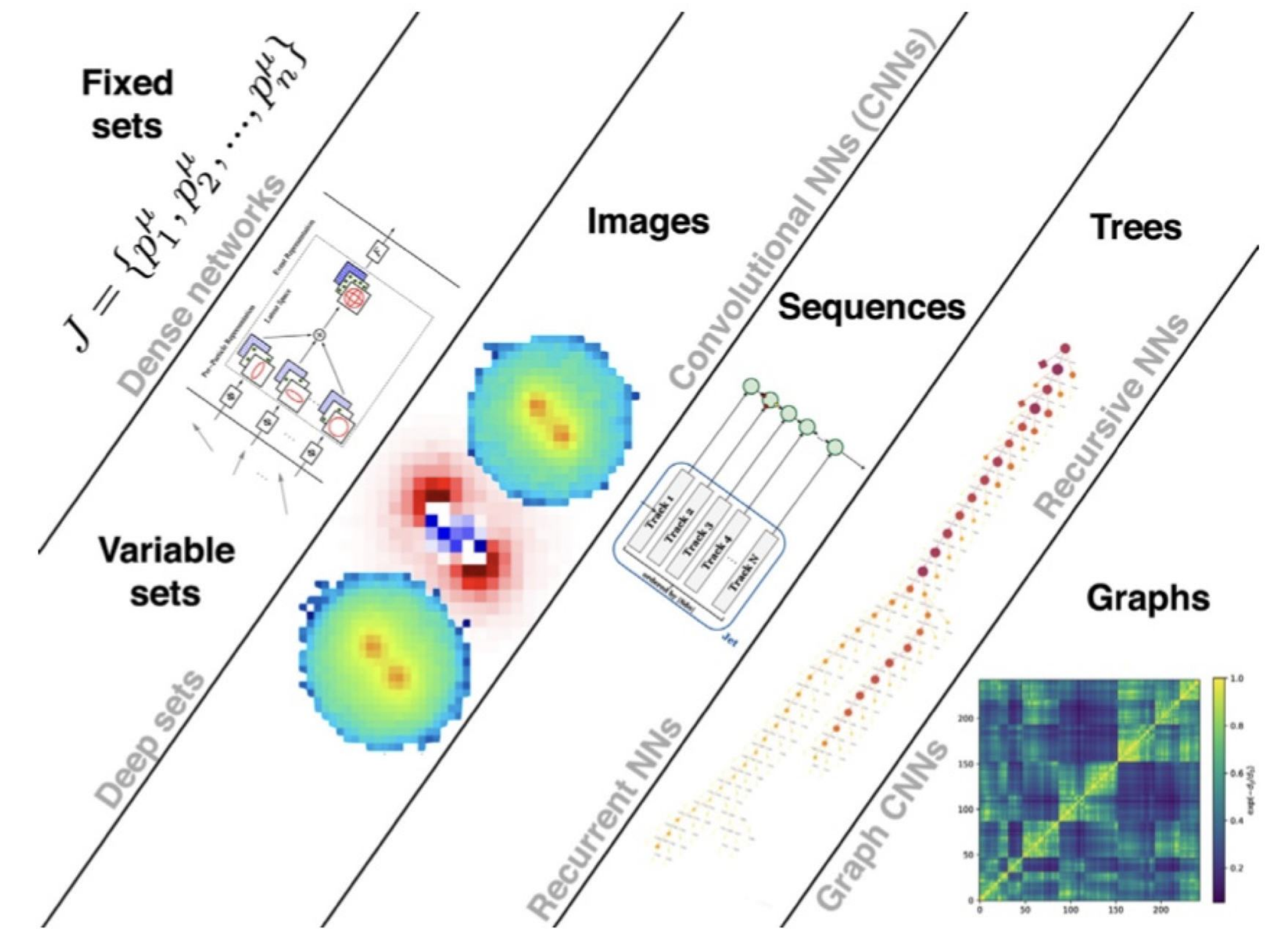
- This alias applies to observables as well.

```
1  n_subjects = parse_observable("subjects.size")
2  n_subjects.read(events).array
```

APIs

Representations

- Set: 1D ($N \times F$)
Use a set of observables to represent an event or a jet.
- Image: 3D ($N \times H \times W \times C$)
Project particles onto a 2D plane.
- Graph: 2D ($N \times P \times F$)
Record particles' features.



[1709.04464] Jet Substructure at the Large Hadron Collider:
A Review of Recent Advances in Theory and Machine Learning

APIs

Representations: Set

```
1 from hml5.approaches import Cut
2 from hml5.events import load_events
3 from hml5.representations import Set
4
5
6 events = load_events("../..//tag_1_delphes_events.root")
7
8 cut1 = Cut("fatjet.size > 0").read(events)
9 cut2 = Cut("ak4fatjet.size > 1").read(events)
10
11 sets = Set("fatjet0.mass", "fatjet0.tau21", "ak4fatjet0,ak4fatjet1.delta_r")
12 sets.read(events, cut1.array & cut2.array)
```

- Names are valid through very part of the workflow.
- Use **read** to extract data from events.

APIs

Representations: Image

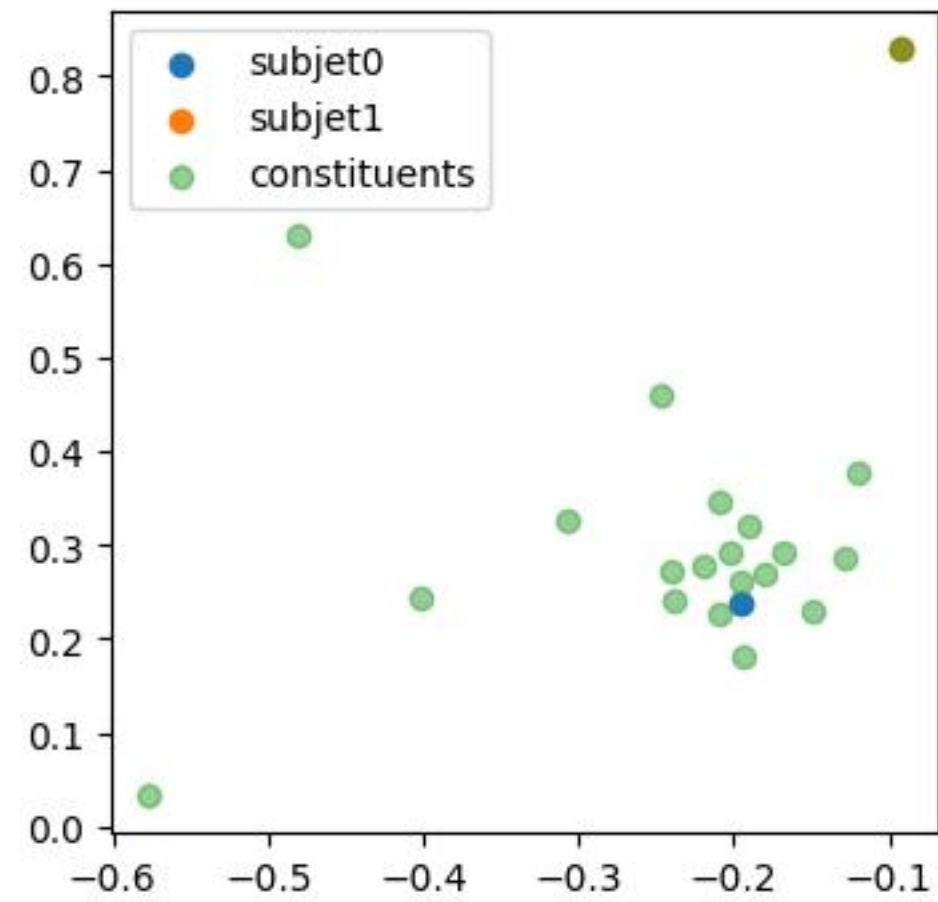
```
1  images = (  
2      Image(  
3          "fatjet0.constituents.eta",  
4          "fatjet0.constituents.phi",  
5          # "fatjet0.constituents.pt",  
6      )  
7      .translate("ak4fatjet0.reclustered0")  
8      .rotate("ak4fatjet0.reclustered1", -90)  
9      .pixelate((32, 32), ((-1.5, 1.5), (-1.5, 1.5)))  
10     .read(events, cut0.array & cut1.array)  
11 )
```

- Recluster the first fatjet and get the first subset.
- Define how images are pre-processed.

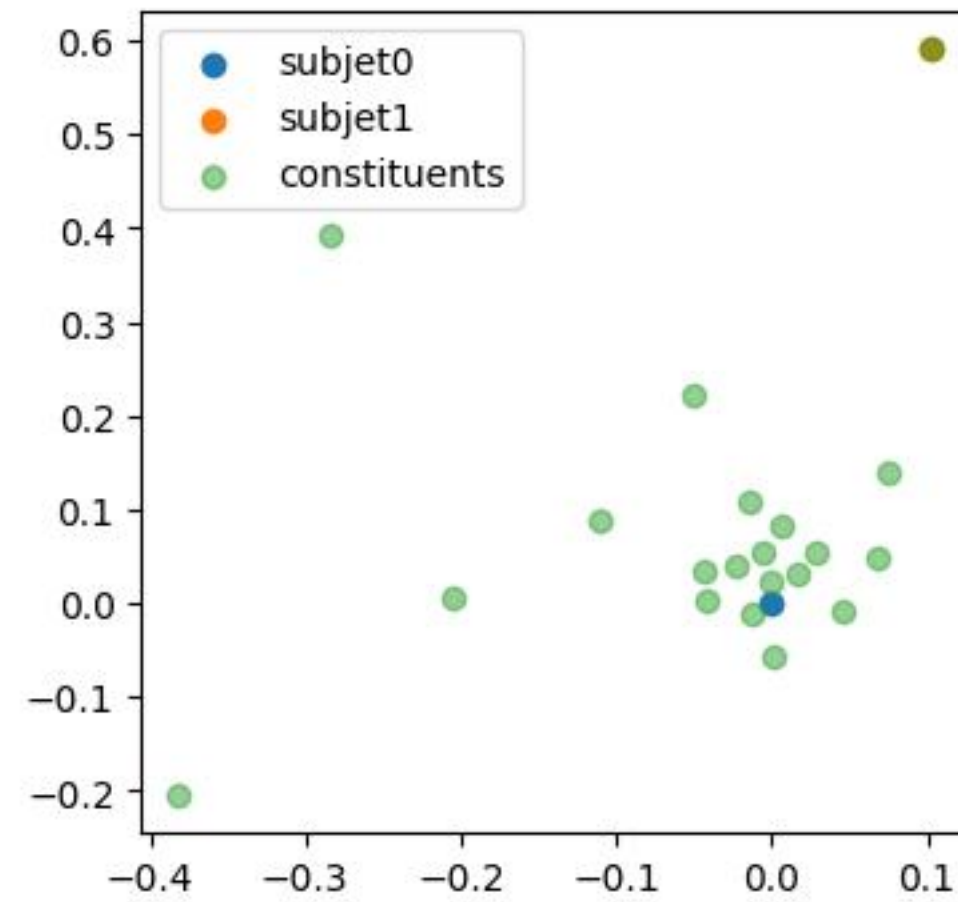
APIs

Representations: Image

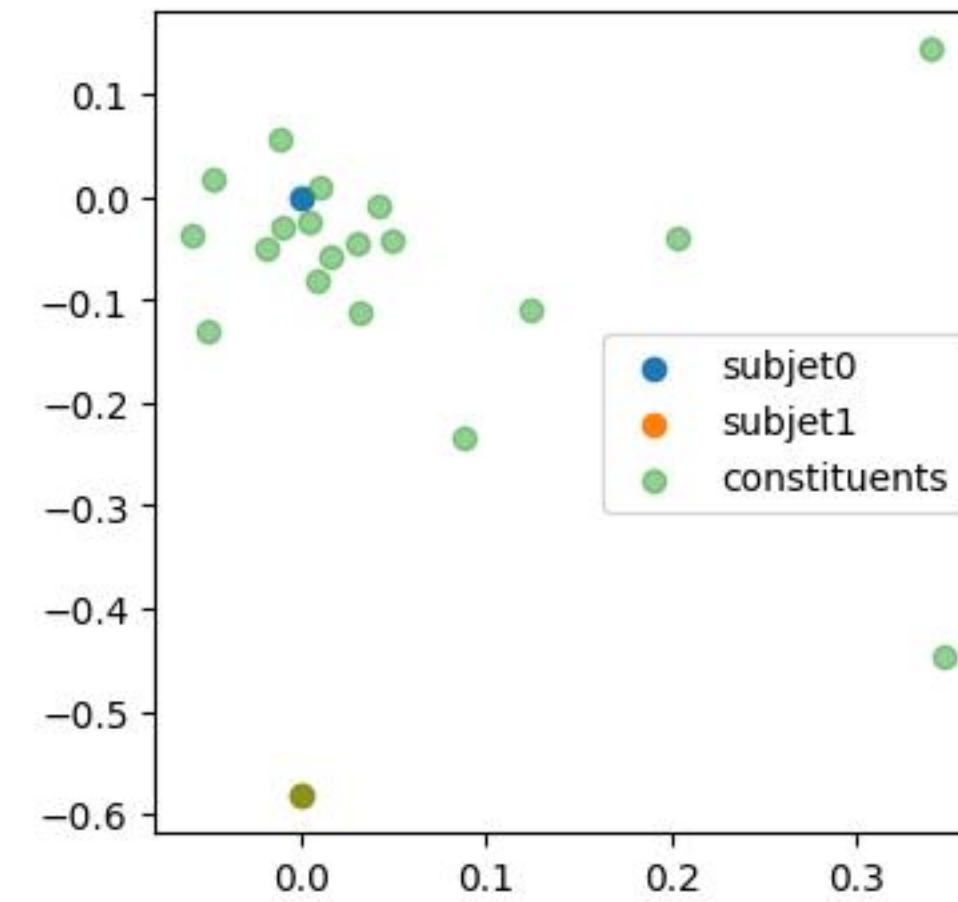
Raw



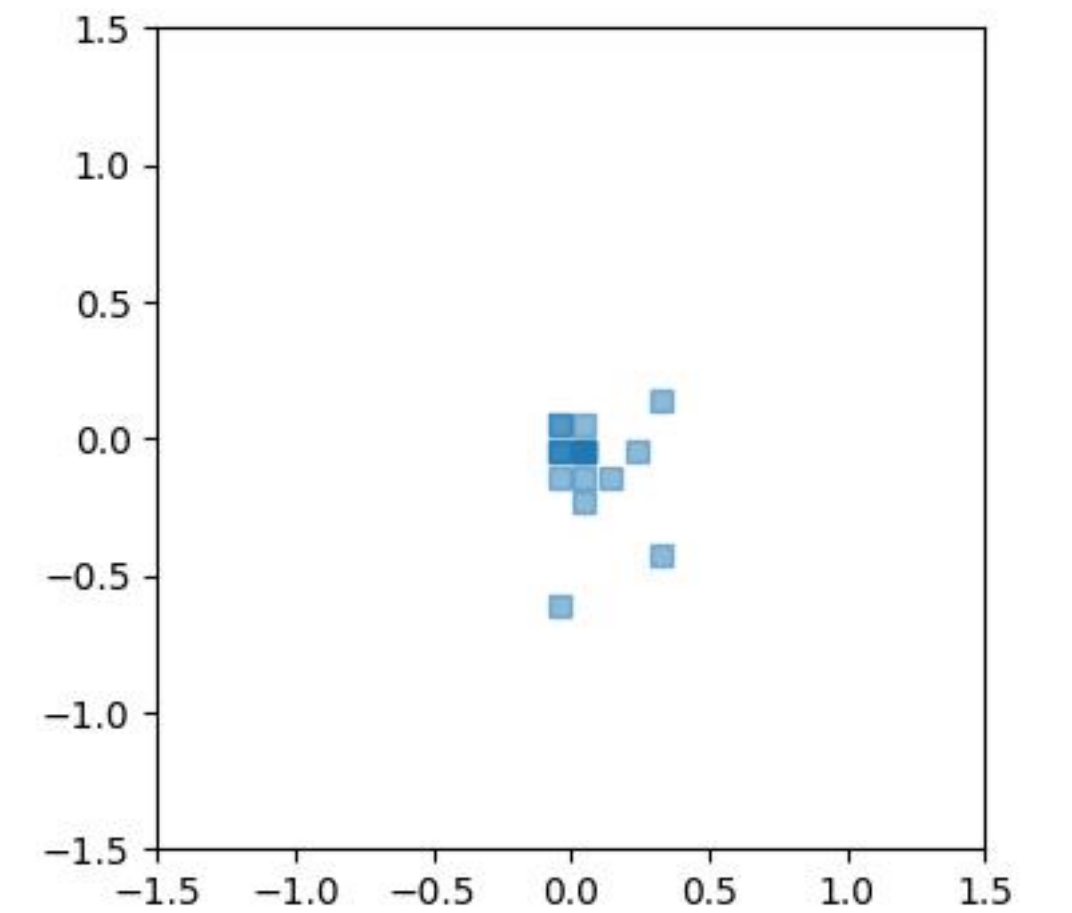
Translated



Rotated



Pixelated

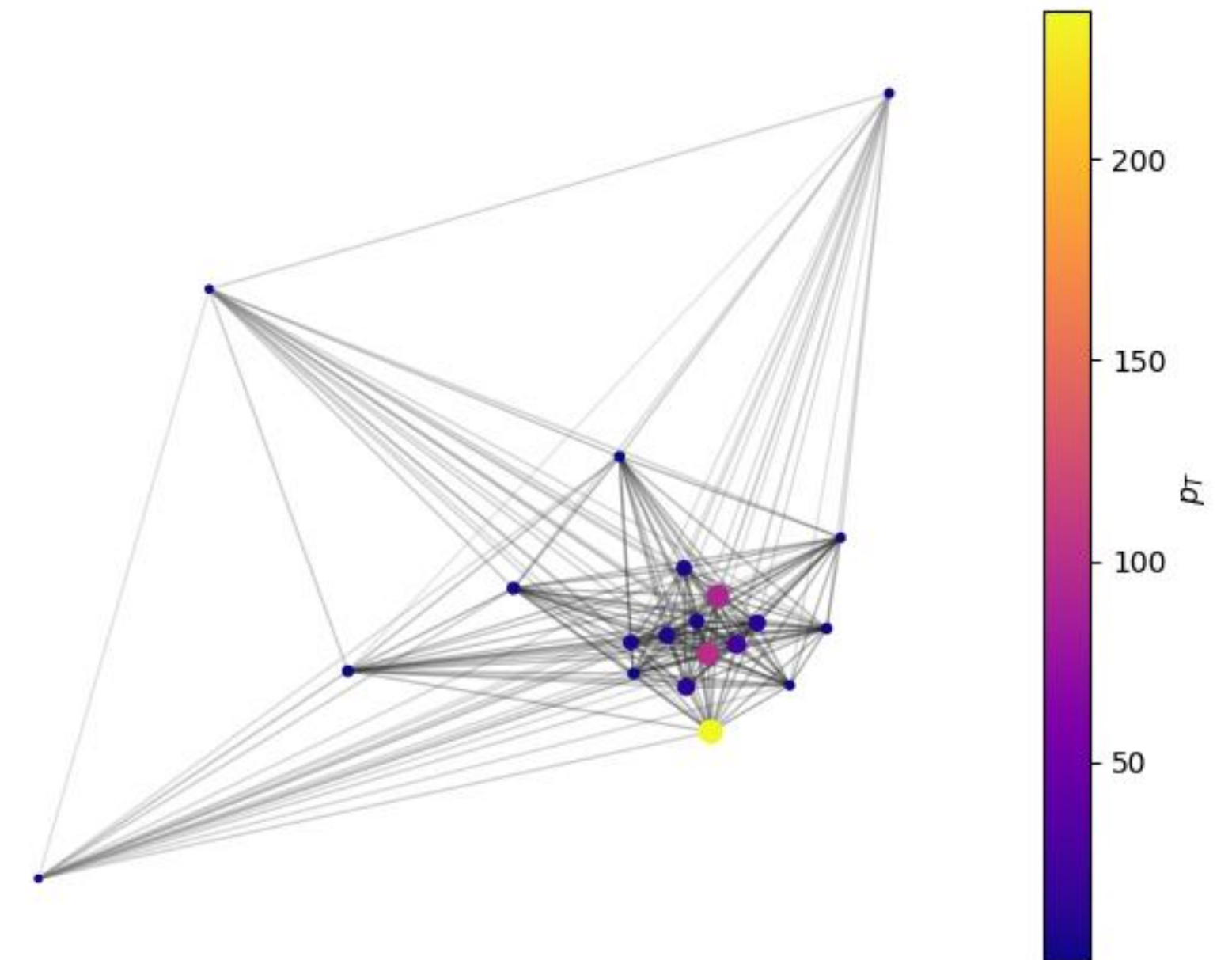


APIs

Representations: Graph

- The same goes for the graph representation.

```
1 graphs = Graph(  
2     Set(  
3         "fatjet0.constituents.eta",  
4         "fatjet0.constituents.phi",  
5         "fatjet0.constituents.pt",  
6     ),  
7     "fatjet0.constituents,fatjet0.constituents.delta_r",  
8 ).read(events, cut.array)
```



APIs

Create datasets

- Create datasets with the data fetched by different representations.

```
1 ds = SetDataset(  
2     [  
3         "fatjet0.mass",  
4         "fatjet0.tau21",  
5         "ak4fatjet0,ak4fatjet1.delta_r",  
6     ]  
7 )  
8 ds.read(events, target=0, mask=cut1.array & cut2.array)  
9 ds.split(0.7, 0.2, 0.1)  
10 ds.save("dataset.ds", piece_size=10)
```

A set dataset:

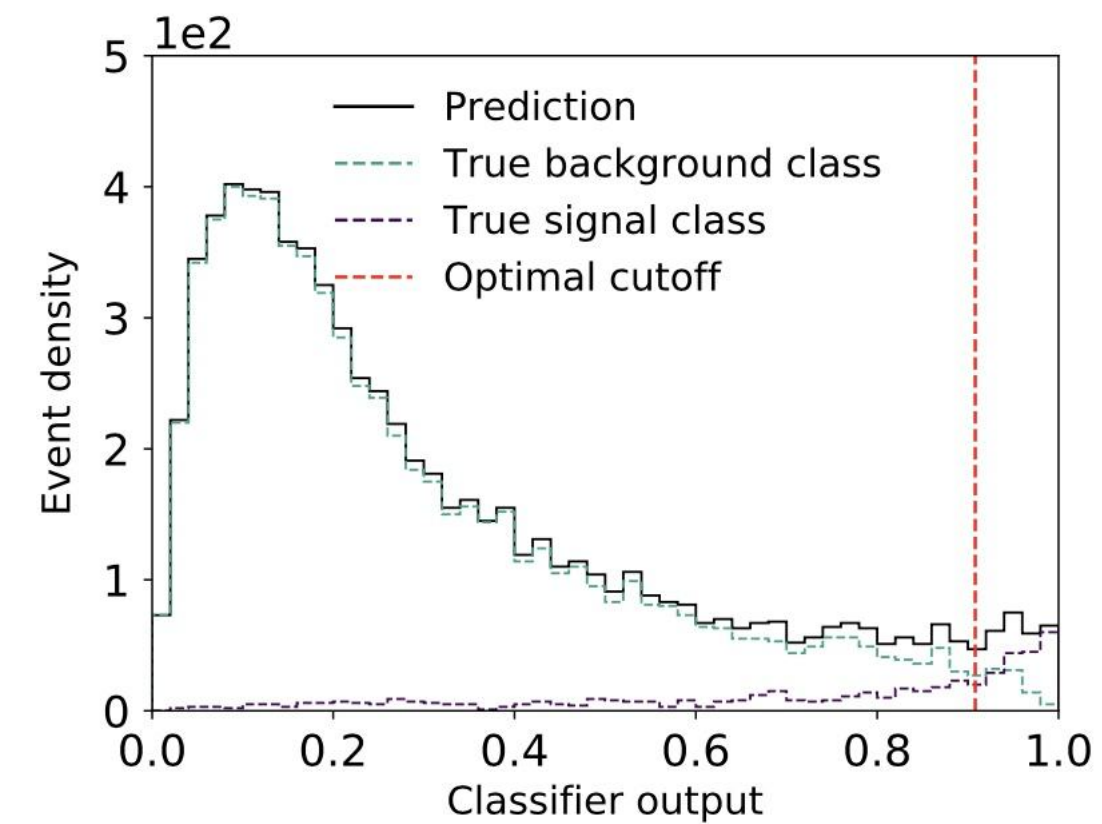
- train: x: 65 * 3 * float32, y: 65 * int32
- test: x: 19 * 3 * float32, y: 19 * int32
- val: x: 10 * 3 * float32, y: 10 * int32

APIs

Approaches

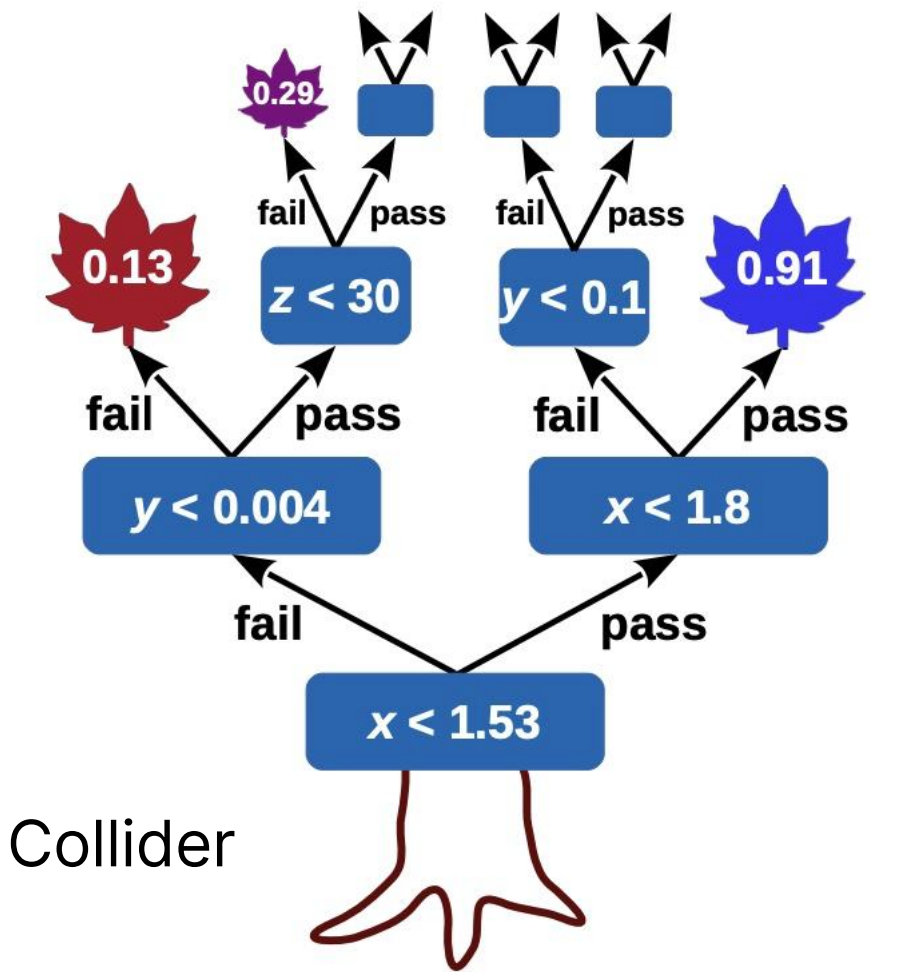
- Designed to contain three kinds of approaches:
 - cut and count
 - tree
 - neural networks

[2108.03125] Beyond Cuts in Small Signal Scenarios

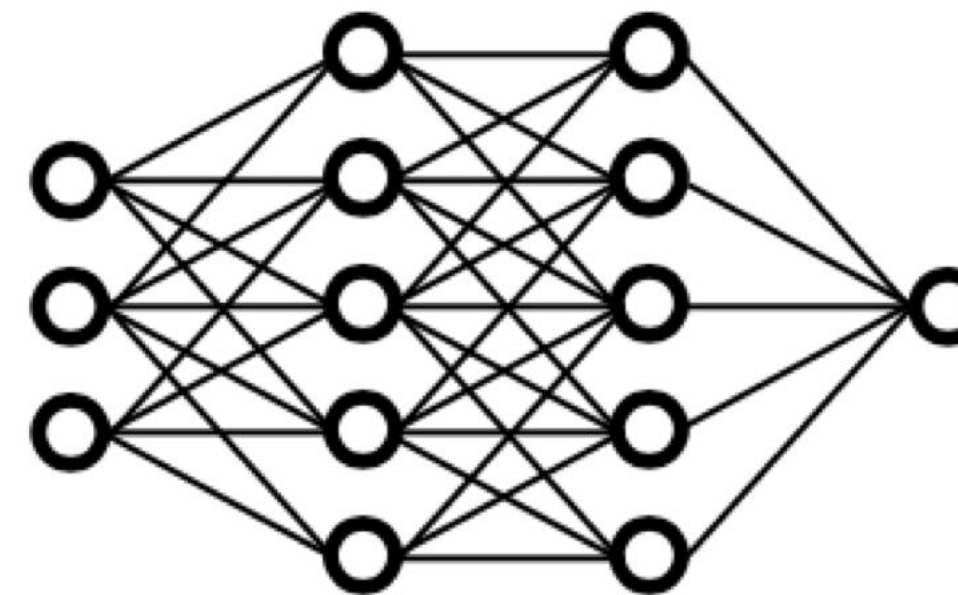


(a) XGBoost with optimized cutoff at 0.9081.

[2206.09645] Boosted decision trees

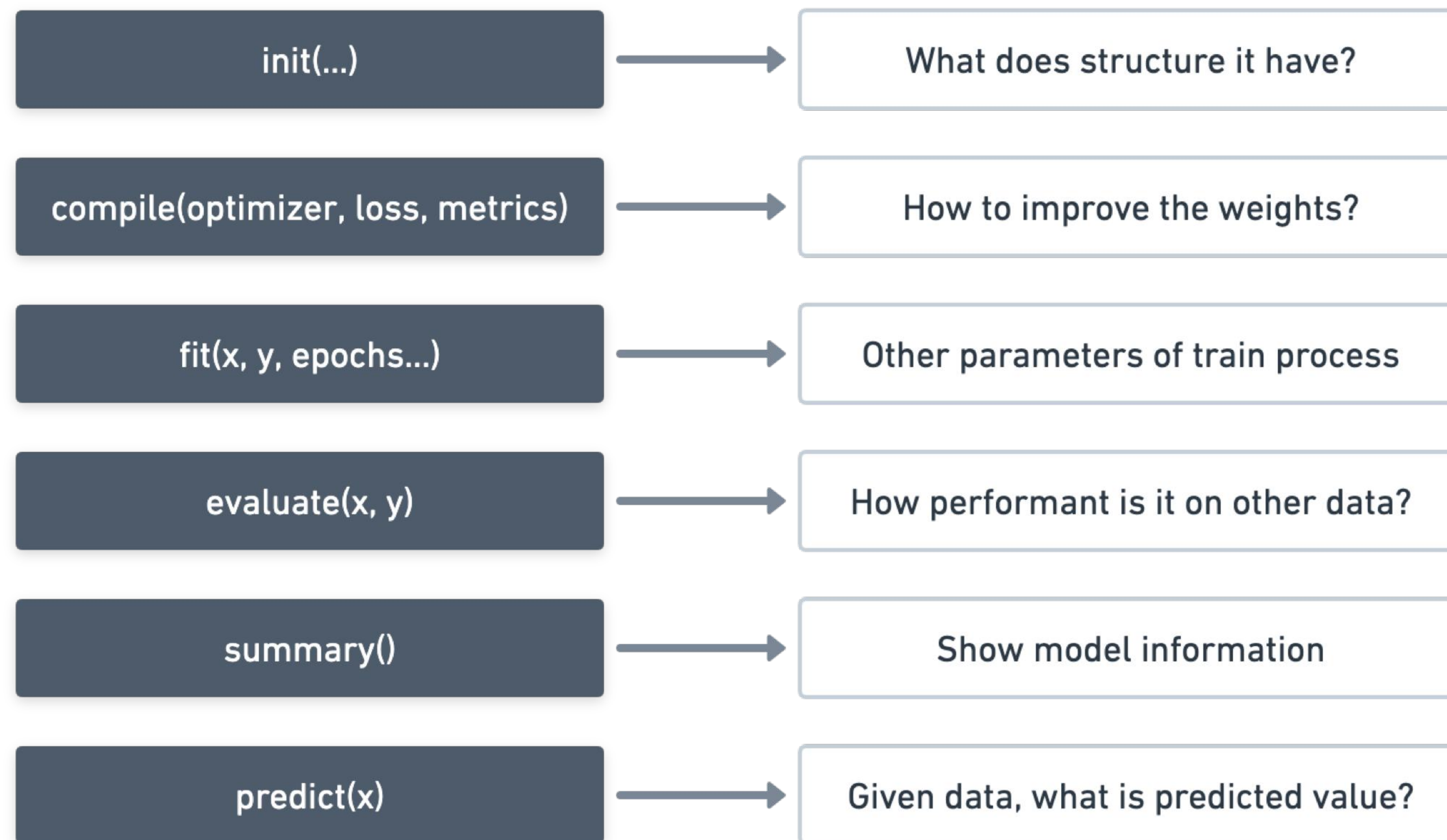


[1709.04464] Jet Substructure at the Large Hadron Collider



APIs

Approaches: protocol



- No need to re-wrap Keras model.

APIs

Approaches: metrics

- **MaxSignificance** calculates the maximum significance under uniform distributed thresholds.

$$\text{significance} = \sqrt{2 \left((S + B) \ln \left(1 + \frac{S}{B} \right) - S \right)}$$

- **RejectionAtEfficiency** calculates the background rejection at a given signal efficiency.

$$1/\varepsilon_b \text{ at } \varepsilon_s = 50\%$$

APIs

Apply approaches

- Initialize methods to define their structure respectively.
- Compile each to determine how to improve itself and monitor performance.
- Fit methods' weights on dataset.

```
1 m1 = BoostedDecisionTree(n_estimators=10)
2 m2 = CutAndCount(n_bins=100)
3 m3 = ToyMLP(input_shape=(x_train.shape[1],))
4
5 m1.compile(
6     loss=CategoricalCrossentropy(),
7     metrics=[
8         CategoricalAccuracy(name="acc"),
9         MaxSignificance(name="max_sig"),
10        RejectionAtEfficiency(name="r50"),
11    ],
12 )
13 m2.compile(...)
14 m3.compile(...)
15
16 m1.fit(x_train, y_train)
17 m2.fit(...)
18 m3.fit(...)
19
```


APIs

Apply approaches

```
1 from tabulate import tabulate
2
3 results1 = method1.evaluate(x_test, y_test)
4 results2 = method2.evaluate(x_test, y_test)
5 results3 = method3.evaluate(x_test, y_test, verbose=2)
6 results = {}
7
8 results["name"] = [method1.name, method2.name, method3.name]
9 for k in results1.keys():
10     results[k] = results1[k] + results2[k] + results3[k]
11
12 print("> Results:")
13 print(tabulate(results, headers="keys", floatfmt=".4f"))
14
```

```
1 > Results:
2 name          loss      acc      max_sig      r50
3 -----
4 boosted_decision_tree  0.2611  0.9586  601.7032  647.3771
5 cut_and_count      4.4163  0.8037  243.9667  33.6241
6 toy_mlp            0.5475  0.9350  111.5401  444.2333
```

- Evaluate methods using metrics defined in compile methods. Could also compile once again to use other metrics.
- Later more metrics will be added to complete benchmark.

Roadmap

v0.4.3(May)

- Prepare basic scaffold.
- Design all modules

v0.5.0(~Aug.)

- Refactor the naming convention.
- New representation: Graph
- New datasets
 - W
 - Z
 - Top
- New models:
 - ParticleNet

v0.6.0(~Sep.)

- Benchmarks
 - Jet tagging
 - New physics

Roadmap

- [2405.02888] HEP ML Lab: An end-to-end framework for applying machine learning into phenomenology studies
- [2303.15920] Probing Heavy Neutrinos at the LHC from Fat-jet using Machine Learning
- <https://github.com/Star9daisy/hep-ml-lab>

```
pip install hep-ml-lab
```

THANK YOU!