CERN Open Data for Machine Learning



Dr. Sercan Şen (Hacettepe University)



NEAR EAST UNIVERSITY February 21, 2020

Sercan.Sen@cern.ch

Outline

- What is machine learning (ML)?
 - Types of ML / classifiers / decision trees
- LHC experiments at CERN
 - Data format
 - High energy physics + ML
- CERN open data
- ML dedicated open data release
 - ML approach to jet identification
 - pixel tracking studies
 - Higgs to bb tagging
- Summary and outlook

What is Machine Learning (ML)?

- Wide field with increasing applications in research and industry
- First ideas already since 1950's
- Machine Learning not new to physics
- In HEP, many 'classic' ML methods are already in use.
 - Boosted decision trees (BDTs) in Higgs search
 - In ALICE, e.g. BDTs for signal extraction for charmed baryons
- Relatively new: Deep learning in physics
 - → Huge progress done in last years
 - → More and more analyses upcoming in all experiments

Development of Machine Learning techniques is incredibly fast

- A lot of progress from tech companies and industry
- Boost from big data, most progress in deep learning
- Many advanced solutions for "human problems" (e.g. image recognition, text understanding) can also be adapted to HEP problems!

Hello World – C++ / Python

Example 1: Hello World Program

```
#include <iostream>
using namespace std;
```

```
int main()
{
    cout << "Hello, World!";
    return 0;</pre>
```

```
}
```

Output

Hello, World!

To Run the code:

g++ example.cpp -o exp ./exp

Source Code

This program prints Hello, world!

print('Hello, world!')

Output

Hello, world!

To Run the code:

python example.py

Hello World – C++ / Python

This program checks whether user input number is positive, negative, or zero and displays the result

```
num = float(input("Enter a number: "))
 if num \geq 0:
      if num == 0:
          print("Zero")
      else:
                                            Rules are written by hand.
          print("Positive number")
 else:
                                            Rule #1:
      print("Negative number")
                                            if the number is greater than zero, then the
                                            number is "positive".
Output 1
 Enter a number: 5
 Positive number
```

Hello World – Machine Learning



Can you write code to tell the difference between an apple and an orange?

Tons of lines of rules are needed to * do that

We need an algorithm that can figure out the rules for us, so we do not have to write them by hand.

 \rightarrow Traine a classifier

def detect_colors(image):
 # lots of code

def detect_edges(image):
 # lots of code

def analyze_shapes(image):
 # lots of code

def guess_texture(image):
 # lots of code

def define_fruit():
 # lots of code

def handle_probability():
 # lots of code

If you change the problem (objects), rules change

What is a Classifier?

It takes some data as input and assigns a label to it as output.



The technique to write the classifier automatically is called supervised learning.

Types of Machine Learning

Supervised Learning algorithms need labeled training data and learn the correct mapping of input data and desired output



Unsupervised Learning tries to find a structure in the data without a priori knowledge of desired outcome



Training Data

Features: Weight and Texture Label: Orange and Apple labeled training data supervised learning

Weight	Texture	Label
150g	Bumpy	Orange
170g	Bumpy	Orange
140g	Smooth	Apple
130g	Smooth	Apple

The more training data you have, the better a classifier you can create

Training Data



To train data, we use a type of classifier called a **decision tree**.



New classifier for a new problem can be created, just by changing the training data!

Much better than writing new rules for each problem.

Decision Trees (I)







The dataset contains a set of 150 records under five attributes - petal length, petal width, sepal length, sepal width and species.



Fisher's Iris Data [hide]

Dataset Order ¢	Sepal length +	Sepal width \$	Petal length +	Petal width +	Species +
1	5.1	3.5	1.4	0.2	I. setosa
2	4.9	3.0	1.4	0.2	I. setosa
3	4.7	3.2	1.3	0.2	I. setosa
4	4.6	3.1	1.5	0.2	I. setosa
5	5.0	3.6	1.4	0.3	I. setosa
6	5.4	3.9	1.7	0.4	I. setosa
7	4.6	3.4	1.4	0.3	I. setosa
8	5.0	3.4	1.5	0.2	I. setosa
9	4.4	2.9	1.4	0.2	I. setosa
10	49	31	15	01	Leetosa

→Three types of flowers

 \rightarrow 50 examples for each type

→150 examples in total

Decision Trees (II)

7.2. Toy datasets

scikit-learn comes with a few small standard datasets that do not require to download any file from some external website.

They can be loaded using the following functions:

<pre>load_boston([return_X_y])</pre>	Load and return the boston house-prices dataset (regression).
<pre>load_iris([return_X_y])</pre>	Load and return the iris dataset (classification).
<pre>load_diabetes([return_X_y])</pre>	Load and return the diabetes dataset (regression).
<pre>load_digits([n_class, return_X_y])</pre>	Load and return the digits dataset (classification).
<pre>load_linnerud([return_X_y])</pre>	Load and return the linnerud dataset (multivariate regression).
<pre>load_wine([return_X_y])</pre>	Load and return the wine dataset (classification).
<pre>load_breast_cancer([return_X_y])</pre>	Load and return the breast cancer wisconsin dataset (classification).

```
- C#2
```

ssen@lxplus700:/afs/cern.ch/woi

```
from sklearn.datasets import load_iris
iris = load_iris()
print iris.feature_names
print iris.target_names
print iris.data[0]
print iris.target[0]
```

import the iris into scikit-learn

https://scikit-learn.org/stable/datasets/ index.html

plus2:tutorials ssen\$ python iris.py ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)'] ['setosa' 'versicolor' 'virginica'] [5.1 3.5 1.4 0.2] 0

- 1. Import dataset
- 2. Train a classifier
- 3. Predict label for new flower
- 4. Visualize the tree

- Examples used to "test" the classifier's accuracy.
- Not part of the training data

Testing Data

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn import tree
```

```
iris = load_iris()
test_idx = [0,50,100]
```

```
#training data
train_target = np.delete(iris.target, test_idx)
train_data = np.delete(iris.data, test_idx, axis=0)
```

```
#testing data
test_target = iris.target[test_idx]
test_data = iris.data[test_idx]
```

```
clf = tree.DecisionTreeClassifier()
clf.fit(train_data, train_target)
```

```
print 'testing data:',(test_target)
print 'predicted label:', (clf.predict(test_data))
```

Output:

plus2:tutorials ssen\$ python iris.py testing data: [0 1 2] predicted label: [0 1 2]

Decision Tree



- How are decision trees built automatically from examples?
- How well do they work in practice?
- The tree asks the questions about the features.
 Choosing good features is one of the important tasks in the first place.

How do we use machine learning in high energy particle physics?

A roadmap for ML in HEP

Machine Learning in High Energy Physics Community White Paper

May 17, 2019

Data science and HEP communities

Abstract: Machine learning has been applied to several problems in particle physics research, beginning with applications to high-level physics analysis in the 1990s and 2000s, followed by an explosion of applications in particle and event identification and reconstruction in the 2010s. In this document we discuss promising future research and development areas for machine learning in particle physics. We detail a roadmap for their implementation, software and hardware resource requirements, collaborative initiatives with the data science community, academia and industry, and training the particle physics community in data science. The main objective of the document is to connect and motivate these areas of research and development with the physics drivers of the High-Luminosity Large Hadron Collider and future neutrino experiments and identify the resource needs for their implementation. Additionally we identify areas where collaboration with external communities will be of great benefit.

Editors: Sergei Gleyzer³⁰, Paul Seyfert¹³, Steven Schramm³²

Contributors: Kim Albertsson¹, Piero Altoe², Dustin Anderson³, John Anderson⁴, Michael Andrews⁵, Juan Pedro Araque Espinosa⁶, Adam Aurisano⁷, Laurent Basara⁸, Adrian Bevan⁹, Wahid Bhimji¹⁰, Daniele Bonacorsi¹¹, Bjorn Burkle¹², Paolo Calafiura¹⁰, Mario Campanelli⁹, Louis Capps², Federico Carminati¹³, Stefano Carrazza¹³, Yi-Fan Chen⁴, Taylor Childers¹⁴, Yann Coadou¹⁵, Elias Coniavitis¹⁶, Kyle Cranmer¹⁷, Claire David¹⁸, Douglas Davis¹⁹, Andrea De Simone²⁰, Javier Duarte²¹, Martin Erdmann²², Jonas Eschle²³, Amir Farbin²⁴, Matthew Feickert²⁵, Nuno Filipe Castro⁶, Conor Fitzpatrick²⁶, Michele Floris¹³, Alessandra Forti²⁷, Jordi Garra-Tico²⁸,

2010: a New Era in Fundamental Science

Exploration of a new energy frontier Proton-proton and Heavy Ion collisions at E_{CM} up to 14 TeV

LHC ring: 27 km circumference

CMS

I O T ER LINC F VIOE DAL

ALICE: Primordial cosmic plasma

ATLAS: Higgs and supersymmetre

CMS: Higgs and supersymmetry

LHCb: Matter-antimatter differen

Highly heterogeneous system Raw data is 100 M channels sampled every 25 ns.



world of physics

Discovery upends 4 JULY 2012 **CERN Press conference**





Nobel Prize in Physics 2013



The Nobel Prize in Physics 2013 was awarded jointly to François Englert and Peter W. Higgs "for the theoretical discovery of a mechanism that contributes to our understanding of the origin of mass of subatomic particles, and which recently was confirmed through the discovery of the predicted fundamental particle, by the ATLAS and CMS experiments at CERN's Large Hadron Collider".

Higgs production and decays at the LHC









Higgs Boson Machine Learning Challenge 2014

~1 million events (rows) 35 features (coloumns)

Eventid	DER_mass_MI	DER_mass_tran	DER_mass_vis	DER_pt_h	DER_deltaet	DER_mass_jet_jet	DER_prod	PRI_jet_subleading_phi	PRI_jet_all_pt	Weight	Label	KaggleSet	KaggleWeight
100000	138.47	51.655	97.827	27.98	0.91	124.711		-2.475	113.497	0.00081448039868	S	t	0.00265331133733
100001	160.937	68.768	103.235	48.146	-999.0	-999.0		-999.0	46.226	0.681041906806	b	t	2.23358448717
100002	-999.0	162.172	125.953	35.635	-999.0	-999.0		-999.0	44.251	0.715742006349	b	t	2.34738894364
100003	143.905	81.417	80.943	0.414	-999.0	-999.0		-999.0	0.0	1.66065435355	b	t	5.44637821192
100004	175.864	16.915	134.805	16.405	-999.0	-999.0		-999.0	0.0	1.90426344118	b	t	6.24533268686
100005	89.744	13.55	59.149	116.344	2.636	284.584		3.106	193.66	0.0254337596084	b	t	0.0834140312717
100006	148 754	28 862	107 782	106.13	0.733	158 359		-2.767	179.877	0.00081448039868	s	t	0.00265331133733
100007	154 916	10.418	94 714	29 169	-000 0	-999.0		-999.0	30.638	0.00572068250088	s	t	0.018636116672
100007	105 504	50.550	100.080	29.109	-333.0	-999.0		-999.0	0.0	1.614803466	b	t	5.29600298518
100008	105.594	50.559	100.989	4.288	-999.0	-999.0		-999.0	167.735	0.000461025356734	s	t	0.00150187015894
100009	128.053	88.941	69.272	193.392	-999.0	-999.0		-999.0	0.0	0.70114133537	b	t	2.29950373735
100010	-999.0	86.24	79.692	27.201	-999.0	-999.0		-2.079	165.64	0.0936590129021	b	t	0.307169523947
100011	114.744	10.286	75.712	30.816	2.563	252.599		-999.0	93.117	0.512739889611	b	t	1.68161144262
100012	145.297	64.234	103.565	106.999	-999.0	-999.0	-	-999.0	0.0	0.665890037882	b	t	2.18389154017
100013	82.488	31.663	64.128	8.232	-999.0	-999.0		-999.0	0.0	0.655921659265	b	t	2.15119866823
100014	-999.0	109.412	14.398	17.323	-999.0	-999.0		-999.0	0.0	0.00572068250088	s	t	0.018636116672
100015	111.026	32.096	75.271	23.067	-999.0	-999.0		-999.0	36.263	0.443597626883	b	t	1.45484847268
100016	114.256	4.351	67.963	47.221	-999.0	-999.0		-999.0	0.0	0.000461281573949	s	t	0.00150270483101
100017	127.861	50.953	77.267	26.967	-999.0	-999.0		-999.0	0.0	1,56163346758	b	t	5.12162357847
100018	-999.0	85.186	68.827	5.042	-999.0	-999.0		-999.0	0.0	1.82316290593	b	t	5.97935067368

Training data

Events are labeled as "signal" (Higgs) or "background" Train a classifier and predict the signal events

Typical proton-proton collision



Add 40 such on top of each other currently. Up to 200 such overlay in the horizon 2025.

Jets: any particle decaying in quark/gluons will result in 2a "jet" of particles in the direction of the original particle.



Point precision ~5 µm to 3mm 100k points 10k tracks / event

North Marian Shit and the

10-100 billion events/year

 $6 \,\mathrm{m}$

image: David Rousseau

LHC / HL-LHC Plan



LHC





Grid CERN

World Wide Web: an information system where documents interlinked by hypertext, and accesible over the internet.

LHC Grid: A global collaboration of computer centres distributes and stores LHC data, giving real-time access to physicists around the world





Running jobs: 268149 Transfer rate: 11.38 GiB/sec





2/15/2013 4:13:20 pm

- D M

Some and Many and State State States



Luch Thep



Event Filtering



Ultra fast decision to keep the relevant data. In hardware and software.

Computing Grid



From Low to High Level Data



The reconstruction of an event goes from the digital signal of the individual sub-detector to a sequence of particles, jets, and high-level features

$\mathsf{RAW} \to \mathsf{RECO} \to \mathsf{AOD} \to \mathsf{NanoAOD}$

image: Jean-Roch Vlimant

Who can access the data?

- opendata.cern.ch launched in November 2014
- LHC collaboration data policies
 - restricted \rightarrow embargo period (~5 years) \rightarrow open
- Over 1.5 Petabytes of open particle physics data
 - datasets, software, VMs, configuration, documentation,...
- Users
 - education: general public, high-school students, masterclasses
 - research: data scientists, physicists

Openly accesible data

Mathematical Reproducible analysis examples

Lightweight and easily readable data format

Good usability for non-physics experts

Developed by CERN-IT and CERN-SIS in collaboration with Experiments

http://opendata.cern.ch
http://github.cernopendata

Learn

Discover the world of open data from particle physics

Visualise

Explore detector events and run basic histogramming

Analyse

Run your own physics analyses, start virtual machines

Welcome to our updated portal CMS Guide to education use of CMS

Open Data

Improving educational content with high school teachers: A field report from our summer students

Glossary

more

CMS Event Display OPERA Event Display CMS Histograms CMS Guide to research use of CMS Open Data ATLAS Higgs Machine Learning Challenge Getting Started with LHCb Open Data Getting Started with ALICE Open Data

more

😸 News 📎

News & Update

2020-02-11 by ATLAS Collaboration

ATLAS Experiment releases 13 TeV Open Data for Science Education News ATLAS 2019-07-18 by CMS Collaboration

CMS releases open data for

Machine Learning

News CMS

News CMS

2018-05-22 by OPERA Collaboration

Release of the first set of data samples by the OPERA Collaboration

2017-12-20 by CMS Collaboration

Observing the Higgs with over one petabyte of new CMS Open Data

First non-LHC data releases

opendata Search CERN

OPERA tau neutrino candidate and multiplicity studies data

Visualize collisions

Click on a name under "Provenance", "Tracking", "ECAL", "HCAL", "Muon", and "Physics" to view contents in table

Interactive histogramming

Dive deeper into the data (I)

Analysis of the di-muon spectrum using data from the CMS detector taken in 2012.

- Rediscover particle resonances in a wide energy range up to the Z boson.
- About 62 M events from data taken at CMS in 2012.
- Only ROOT as dependency.
- NanoAOD format.
- Analysis code in Python, C++ or a Jupyter notebook.

Dive deeper into the data (II)

opendata

Higgs-to-four-lepton analysis example using 2011-2012 data

Jomhari, Nur Zulaiha; Geiser, Achim; Bin Anuar, Afiq Aizuddin;

Cite as: Jomhari, Nur Zulaiha; Geiser, Achim; Bin Anuar, Afiq Aizuddin; (2017). Higgs-to-four-lepton analysis example using 2011-2012 data. CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.JKB8.RR42

Q

Software Analysis CMS Accelerator CERN-LHC

Description

Search

This research level example is a strongly simplified reimplementation of parts of the original CMS Higgs to four lepton analysis published in Phys.Lett. B716 (2012) 30-61, arXiv:1207.7235.

The published reference plot which is being approximated in this example is https://inspirehep.net/record/1124338/files/H4I_mass_3.png. Other Higgs final states (e.g. Higgs to two photons), which were also part of the same CMS paper and strongly contributed to the Higgs boson discovery, are not covered by this example.

The example consists of different levels of complexity. The highest leve minimal understanding of the content of this paper and of the meaning educational exercises. The lower levels might also be interesting for ed with the linux operating system and the ROOT analysis tool.

Use with

The example uses legacy versions of the original CMS datasets in the At publication due to improved calibrations. It also uses legacy versions o but not identical to, the ones in the original publication. These legacy d in many later CMS publications.

/DoubleElectron/Run2011A-12Oct2013-v1/AOD

/DoubleMu/Run2011A-12Oct2013-v1/AOD

Data Format

SW dependencies?

Virtual machines

	CMS-Open-Data-1.2.0 [Running]		
Applications Menu Terminal - cms-opendata			09:11 💽
Terminal - cms-opendata@lo	ocalhost:~/ssen	000	
<u>File Edit View Terminal Go H</u> elp			
[cms-opendata@localhost ssen]\$ root		<u>^</u>	
Welcome to R00⊤ 6.10/09	http://root.cern.ch		
(c) 1995	-2017 The ROOT Team		
Built for linuxx8664gcc /// States in the second se			
Try '.help', '.demo', '.lisense', '.ch	edits', '.quit'/'.q'		
root [0]		$\sim \pm$	
To the second		Oracle VM VirtualBox Manager	
			ð
	New Settings Discard	Show	Machine Tools
	CMS-Open-Dat	Welcome to VirtualBox!	
Conception: Re	ene Brun, Fons Rademake	The left part of this window lists	00
Core Engineeri Lorenzo Mon	ng: Rene Brun, Fons Radi neta, Vassil Vassilev, Gera	machine groups on your	
Wouter Verke Paul Russo, A	erke, Timur Pocheptsov, M Andrei Gheata, Anirudha E	computer.	
		The right part of this window represents a set of tools which	
		are currently opened (or can be	
		machine. For a list of currently	
		available tools check the corresponding menu at the right	
		side of the main tool bar located	
		list will be extended with new	
		tools in future releases.	

Software tutorials: ROOT

	ROOT Object Browser	000
rowser Eile Edit View Opt	ions <u>T</u> ools	Hel
lles	Canvas 1 🙁 Editor 1 👘	
Draw Option:		
Iroot		
PROOF Sessions		
ROOT Files		
aroot:/Acepublic.cern.ch/Acos/ope	ndataka	
T MetaData,1		
ParameterSets;1		
P Brendage 1	>	
P Events 1		
EventAustory		
EventSalactor		
Branchi istindenes		
All 1 Global Priocer Ibject M	an Record	
Aledm Trigger Results Trigg	er Besud	
Altrigger Trigger Event _htt	riggerS	
EBDigi Collection_select	Digi_sek	
EEDigi Collection_select	Digi_14	· ·
- HoalNoise Sury hary_hoal	noiseCommand	
L1 GlobalTrigeer Readout	Record_+1 Command (local):	*
A second "		

- Very positive feedback from hands-on in ROOT workshops
 - Teaching usage of software with real physics results

```
histo = ROOT.RDataFrame("Events", "Run2012BC_DoubleMuParked_Muons.root")
    .Filter("nMuon == 2")
    .Filter("Muon_charge[0] != Muon_charge[1]")
    .Define("Dimuon_mass",
            "InvariantMass(Muon_pt, Muon_eta, Muon_phi, Muon_mass)")
    .Histo1D("Dimuon_mass")
```

ROOT tutorials: Di-muon, Higgs to four leptons, introductory co46se

CMS Open Data in use for ROOT tutorials

 Providing examples very close to the user's use-cases

Data Format

AOD vs NanoAOD

- **AOD:** Format of primary samples provided on the portal
 - Serialized C++ objects
 - Only fully accessible using CMS software and ROOT
 - Large files with much information (~500 kB/event)
 - Powerful but complex

• NanoAOD: Reduced format often used in recent CMS analyses

- Basic types (floats, integers, ...) or arrays thereof
- Accessible with any library capable to read ROOT files
- Smaller files with reduced information (~1-2 kB/event)
- Used by actual CMS analyses and easily accessible for everyone

• Tool to convert Run 1 AOD files to subset of NanoAOD format for education and outreach

• Ongoing effort to provide Run 1 data in NanoAOD format for research

/ariable	Туре	Description
nMuon	unsigned int	Number of muons in this event
Muon_pt	float[nMuon]	Transverse momentum of the muons (stored as an array of size nMuon)
Muon_eta	float[nMuon]	Pseudorapidity of the muons
Muon_phi	float[nMuon]	Azimuth of the muons
Muon_mass	float[nMuon]	Mass of the muons
Muon charge	int[nMuon]	Charge of the muons (either 1 or -1)

Muon collection in NanoAOD format

Jupyter notebooks

CMS education activities using notebooks and CMS open data

Run analyses on the clouds

Physics publications in peer-reviewed journals

2010 vielding a sample of 768 687 events containing a high-quality central jet with transverse momentum

First analyses by theorists (Jesse Thaler *et_al*, MIT)

Highest-energy particle-collision data ever released through open access.

The ATLAS Collaboration makes public 10 inverse femtobarns (fb⁻¹) of the **13 TeV** data.

Corresponds to about 1 quadrillion proton-proton collisions (that's 1 followed by 15 zeros), or **500 thousand produced Higgs bosons**. Approximately the same amount of data that the ATLAS Collaboration used to discover the Higgs boson in 2012.

Open Data ML release

CMS releases open data for Machine Learning

2019-07-18 by CMS Collaboration

News

The CMS Collaboration at CERN is happy to announce the release of its fourth batch of open data to the public. With this release, which brings the volume of its open data to more than 2 PB (or two million GB), CMS has now provided open access to 100% of its research data recorded in proton-proton collisions in 2010, in line with the collaboration's data-release policy. The release also includes several new data and simulation samples. The new release builds upon and expands the scope of the successful use of CMS open data in research and in education.

http://opendata.cern.ch/docs/cms-releases-open-data-for-machine-learning

In this release, CMS open data address the ever-growing application of machine learning (ML) to challenges in highenergy physics. According to a recent paper, collaboration with the data-science and ML community is considered a high-priority to help advance the application of state-of-theart algorithms in particle physics.

ML studies dedicated datasets on CERN Open Data

opendata Q Search CERN CMS x datascience x Dataset x Sample with jet properties for jet-flavor and other jet-related ML studies JetNTuple_QCD_RunII_13TeV_MC include on-demand datasets Four derived The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of 13 TeV generated with Pythia 8. The particles emerging from the Filter by type collision... V Z Dataset 4 datasets from Derived 4 Dataset Derived CMS Software 4 Tool 4 official 2016 CMS Samples with full event information including tracker hits for tracking, ML, and top . . Iter by experiment quark tagging studies simulation ATLAS 1 Samples in this record are in a custom root ntuple format and contain the position of the hits and 4 CMS information from the generator-level objects associated to the tracker hits. The samples can be US.... (ROOT & HDF5) ilter by year Dataset Derived CMS 2019 4 Jet flavor Filter by file type Sample with tracker hit information for tracking algorithm ML studies h5 3 TTbar_13TeV_PU50_PixelSeeds root з studies The dataset consists of a collection of pixel doublet seeds, i.e. the hit pairs that could belong to the same particle. The compatibility between two hits is evaluated only on the basis of Filter by keywords geometri... Top tagging datascience 4 Dataset Derived CMS Pixel tracking Sample with jet, track and secondary vertex properties for Hbb tagging ML studies studies HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC

The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of 13 TeV generated with Pythia 8. It has been produced for developing machi...

Dataset Derived CMS

http://opendata.cern.ch/record/12102

About -

• H(bb) tagging

jet-related ML studies

Sample with jet properties for jet-flavor and other jet-related ML studies JetNTuple_QCD_RunII_13TeV_MC

Description: The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of 13 TeV generated with Pythia 8.

Dataset characteristics 22554294 entries. 244 files. 204.6 GB in total.					
Dataset semantics					
Variable	Туре	Description			
jetPt	Float_t	Transverse momentum of the jet.			
jetEta	Float_t	Pseudorapidity (ŋ) of the jet.			
jetPhi	Float_t	Azimuthal angle (ϕ) of the jet.			
jetMass	Float_t	Mass of the jet.			
jetGirth	Float_t	Girth of the jet (as defined in arXiv:1106.3076 [hep-ph]).			
jetArea	Float_t	Catchment area of the jet; used for jet energy corrections.			
jetRawPt	Float_t	Transverse momentum of the jet before the energy corrections.			
jetRawMass	Float_t	Mass of the jet before the energy corrections.			
jetLooseID	UInt_t	Binary variable indicating whether the jet passes 'loose' criteria f			
ietTightID	I lint t	Rinary variable indicating whether the jet passes 'tight' criteria fo			

 Training data produced with JetNtupleProducerTool

ML approach to jet identification

ML approach to jet identification

- Re-train ResNet-50 to identify the origin of jets
- Inputs are jet images = pixelated
 versions of calorimeter hits in 2D (η, Φ)

Note: averaged over 10k jets; 1 jet gives a *sparse* image

Higgs to bb

opendata CERN Search Q

Sample with jet, track and secondary vertex properties for Hbb tagging ML studies HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC

Dataset Derived Datascience CMS CERN-LHC

Parent Dataset: /BulkGravTohhTohbbhbb_narrow_M-600_13TeV-madgraph/RunIISummer16MinIAODv2-PUMoriond17_80X_mcRun2_asymptotic_2016_TranchelV_v6_ext1-v1/MINIAODSIM

Description

The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of 13 TeV generated with Pythia 8. It has been produced for developing machine-learning algorithms to differentiate jets originating from a Higgs boson decaying to a bottom quark-antiquark pair (Hbb) from quark or gluon jets originating from quantum chromodynamic (QCD) multijet production.

The reconstructed jets are clustered using the anti-kT algorithm with R=0.8 from particle flow (PF) candidates (AK8 jets). The standard L1+L2+L3+residual jet energy corrections are applied to the jets and pileup contamination is mitigated using the charged hadron subtraction (CHS) algorithm. Features of the AK8 jets with transverse momentum pT > 200 GeV and pseudorapidity $|\eta| < 2.4$ are provided. Selected features of inclusive (both charged and neutral) PF candidates with pT > 0.95 GeV associated to the AK8 jet are provided. Additional features of charged PF candidates (formed primarily by a charged particle track) with pT > 0.95 GeV associated to the AK8 jet are also provided. Finally, additional features of reconstructed secondary vertices (SVs) associated to the AK8 jet (within $\Delta R < 0.8$) are also provided.

Derived datasets (ROOT & HDF5):

http://opendata-dev.web.cern.ch/record/12102

- 182 files, 245 GB, 18 million total entries (jets)
- event features, e.g. MET, ρ (average density)
- big jet features, e.g. mass, p_T, N-subjettiness variables
- ▶ particle candidate features, e.g. p_T , η, φ (*for up to 100 particles*)

charged particle / track features, e.g. impact parameter (for up to 60 tracks) secondary vertex features, e.g. flight distance (for up to 5 stertices)

Higgs to bb

Dataset semantics

Variable	Туре	Description
event_no	UInt_t	Event number
npv	Float_t	Number of reconstructed primary vertices (PVs)
ntrueInt	Float_t	True mean number of the poisson distribution for this event from which the number of interactions in each bunch crossing has been sampled
rho	Float_t	Median density (in GeV/A) of pile-up contamination per event; computed from all PF candidates of the event
sample_isQCD	Float_t	Boolean that is 1 if the simulated sample corresponds to QCD multijet production
fj_doubleb	Float_t	Double-b tagging discriminant based on a boosted decision tree calculated for the AK8 jet (see CMS-BTV-16-002)
fj_eta	Float_t	Pseudorapidity η of the AK8 jet
fj_gen_eta	Float_t	Pseudorapidity η of the generator-level, matched heavy particle: H, W, Z, top, etc. (default = -999)
fj_gen_pt	Float_t	Transverse momentum of the generator-level, geometrically matched heavy particle: H, W, Z, t, etc. (default = -999)
fj_isBB	Int_t	Boolean that is 1 if two or more b hadrons are clustered within the AK8 jet (see SWGuideBTagMCTools)

links to variable definitions / methods

b-tagging

b hadrons have long lifetimes: travel O(mm) before decay!

- displaced tracks
- secondary vertices
- soft leptons

Combined taggers

b-tagging

b hadrons have long lifetimes: travel O(mm) before decay!

- displaced tracks
- secondary vertices
- soft leptons

Super Combined taggers

Deep bb tagging

Large performance gain over previous algorithm (BDT)

Summary and Outlook

- Data science and particle physics communities closely work together.
- A white paper has been published recently.
 - Future research and development areas for machine learning in particle physics.
 - Roadmap for software, hardware resource requirements.
- CERN open data released new data dedicated to Machine Learning studies.