Machine Learning within (DESY) CMS

A general overview with special emphasis on local activities



Dirk Krücker - DESY CMS group 16.9.2019





Overview

CMS activities

CMS is very active in applying Machine Learning and especially Deep Learning to a variety of subjects

- About ~25 dedicated ML contributions on conferences in 2019-present (not including physics analyses)
- Annual dedicated CMS Machine Learning workshops
- ML is transforming the field and the traditional workflows
 - New tools and new terminology

Main fields of applications

- **Object tagging** and object calibration especially Jets and taus are **flagship** applications
- Physics analyses
- In addition
 - Triggers
 - Reconstruction
 - Data Quality Monitoring
 - Simulation
 - Computing work flows

Jet Tagging

DeepJet

The latest approach on b-quark jet tagging

CMS uses a Particle Flow approach

- Identify each particle by the combined information of all sub-detector
- combine to jets (anti-kt)

Jet identification had been one of the first application of Deep Learning in HEP

- **b-quark jets** vs g/u/d/s-jets etc.
- The classical approach looks for secondary vertices
- There is plenty of information in the correlation between the individual tracks and vertices and global event characteristics





Evolution of CMS b-tagging

A short history (for more details M. Verzetti ML4Jets WS 2018)

- 2015 handcrafted **CSVv2(nn)** and cMVAv2(BDT) •
 - Mainly the same information
- 2017 deepCSV •
- mid 2017 **deepJet** (DeepFlavour) •
 - The gain is due to the usage of all available particle information in a deep network





deepJet

Neutral (8 features) x25

see next page

b

С

g



Training DeepJet

Big data

• Training needs large amount of data

- Different weight initialization
 - Training is stable and reproducible

Scale factors

Back to earth

- Training on MC
- Performance on data
 Scale Factors (SF) that describes the difference between simulation and data
 - Good performance on data but smaller efficiency

Top etc. tagging

Highly-boosted top quarks

- A top decays into a b quark and a W \rightarrow lep,had
- If the momentum is large the form a common jet
 - Boosted topology in fat jets (AK8)
- Top-tagging \Leftrightarrow identify such jets

Typical problem how to handle a varying number of particles, here: 1d CNN

DeepAK8

- Up to 100 particles
- Each particle comes with 42 features
- 14 layers is indeed deep
- Up to 7 SV with 7 features
- Trained with 50 million jets
- Residual net (short cuts between layers)

CMS publication will appear soon

mxnet

Convolution allows to handle a varying number of particle within a jet

How to feed in a changing number of objects

A common subject

Many Deep Learning approaches have been developed for imaging, i.e. dense regular grid of points; or text processing, i.e. sequences

- Not straightforward how to connect to typical HEP problems
 - Jet images in (η, ϕ) plane are sparse
 - Varying number of particles in cone from jet to jet
- Standard solution so far
 - Recurrent nets: p_T ordered particles
 - Convolution: particles as 1d string of inputs
- New approach: Geometric learning or

Graph networks

(Hot subject on this year NIPS dozens of papers)

- H. Qu and L. Gouskos. "ParticleNet: Jet Tagging via Particle Clouds" https://arxiv.org/abs/1902.08570
 - Bsed on Y. Wang et al. "Dynamic Graph CNN for Learning on Point Clouds"
- M. Fey and J.E. Lenssen. "Fast Graph Representation Learning with PyTorch Geometric" <u>https://github.com/rusty1s/pytorch_geometric</u>

EdgeConv network

Particle Cloud

- Compute k-nn Graph •
 - Defines a local neighborhood on which • a convolution operation work
 - For the first layer this is the $(\Delta \eta, \Delta \phi, p_T)$ plane wrt. jet axis $e'_{ijm} = \text{ReLU}(\boldsymbol{\theta}_m \cdot (\mathbf{x}_j - \mathbf{x}_i) + \boldsymbol{\phi}_m \cdot \mathbf{x}_i)$
- Apply EdgeConv, here •

$$x'_{im} = \max_{j:(i,j)\in\mathcal{E}} e'_{ijm}$$

Dynamic Graph CNN • means the recalculation of the k-nn in the x' space

EdgeConv network

Thanks to Leonid Didukh

Tau-lepton tagging

A newly started project with Graph Networks – no official CMS results yet

- CMS investigates deep NN for tau-lepton identification
- Similar complexity as the shown b tagger
- We started to investigate a Graph Network approach
- Promising results similar performance with a model that contains 10 times less parameters
 - Time advantage in interference

Searches

Introduction | D. Krücker

Searches

Examples for ML methods in searches

Two example analyses (there are many more using ML) PI

- Supersymmetry
- Higgs

Please note!

- These are ongoing analyses
- Results are not approved
- Therefore no numbers, no data plots

Ex.1: Search for supersymmetry in events with one lepton

Search for supersymmetry in events with one lepton and multiple jets exploiting the angular correlation between the lepton and the missing transverse momentum in proton-proton collisions at \sqrt{s} = 13 TeV

- SUSY signal: T1tttt
 2 gluinos that decay to 4 top quarks
 - + 2 neutralinos

 Important SM background top-anti-top pair (ttbar) with 2 leptons where

one lepton is missed in the reconstruction

published Run II result on 2016 data

SUSY deltaPhi

Previous approach

- Can we do better with a neural network?
 - Sure but
- How to connect the NN response to the observed data to estimate the background in the signal region?

Estimate the background in the high dPhi region by looking into the low dPhi region at different jet multiplicities

Multiclassification

Data augmented background estimation

- 4 classes
 - 3 class represent different backgrounds
 - 1 class is the search region
 - The simulation is normalized to the data in the background classes by solving the equation system

 $Y_i = A_{ij}X_i$ $Y_i \rightarrow \text{ data counts}$ $A_{ij} \rightarrow \text{MC counts}$ $X_i \rightarrow \text{ scale factor for each background}$ $i \rightarrow \text{ number of CRs}$ $j \rightarrow \text{ background number}$

background	α	β	γ
1500_1000	0.84 ± 0.01	$1.0\pm~0.05$	0.72 ± 0.03
1500_1200	0.84 ± 0.01	0.98 ± 0.04	0.72 ± 0.03
1600_1100	0.84 ± 0.01	0.98 ± 0.05	0.73 ± 0.03
$1700_{-}1200$	0.84 ± 0.01	0.98 ± 0.05	0.73 ± 0.03
1800_1300	0.84 ± 0.01	0.98 ± 0.04	0.71 ± 0.03
1900_100	0.84 ± 0.01	0.97 ± 0.05	0.71 ± 0.03
1900_1000	0.84 ± 0.01	0.96 ± 0.04	0.71 ± 0.03
1900_800	0.84 ± 0.01	0.94 ± 0.05	0.72 ± 0.03
2200_100	0.84 ± 0.01	0.96 ± 0.05	0.73 ± 0.03
2200_800	0.84 ± 0.01	0.96 ± 0.04	0.71 ± 0.03

- Creates a nice data/MC agreement -> estimate Background in signal region (not shown here)
- Independent of the signal point

How to know what is relevant?

Opening the neural net black box in a Higgs search

- Complex neural nets do not tell us how they come to a decision
 - 10th thousands of parameters
 - Several dozens of input variables
- Important to understand what is driving the decision
 - If the network is just a chain of matrix multiplication and function mappings why not just do a Taylor expansion

• arXiv:1512.02479

Deep Taylor Decomposition

- Higgs $\rightarrow \tau \tau$
 - Observation published last year (CMS) *Phys. Lett. B* 779 (2018) 283
 - New study with a multiclass NN ongoing

Relevance propagation of variables as applied in $H\to\tau\tau$

Teresa Lenz, Mareike Meyer, Alexei Raspereza for the HIG-18-032 team

Analysis strategy

- four most sensitive final states of $\tau\tau$ -pair studied: eµ, e τ_h , µ τ_h , $\tau_h\tau_h$
- loose baseline selection (trigger requirements, suppression of large backgrounds)
- multi-class NN with 2 signal classes (ggH, qqH) & several background classes (control regions)
- selection & validation of NN input variables based on 1D and 2D GoFs relevance propagation

cross section

 V_{τ}

W

e, μ , d

 \overline{V}_{e} , \overline{V}_{u} , \overline{U}

Relevance propagation: idea

S.Wunsch, R. Friese, R.Wolf, and G. Quast, Computing and Software for Big Science 2 (Sep, 2018) 5, doi:10.1007/ s41781-018-0012-, arXiv:1803.08782

- relate output space of NN to input space
- identify characteristics of input space that have large influence on output for a given task
- decompose NN function into Taylor expansion in each element of the input space
- Taylor coefficients contain information about the sensitivity of the NN response to the inputs
- dependence on phase space: mean of absolute values of Taylor coefficients evaluated for all elements of test sample

set of input variables, evaluated for element k of test sample

Size of the Taylor coefficients as sensitivity metric

- first order Taylor coefficient: influence of single input elements
- second order Taylor coefficent: influence of pair-wise or auto-correlations

$$T(x,y) = f(a,b) + (x-a)f_x(a,b) + (y-b)f_y(a,b) + rac{1}{2!}\Big((x-a)^2f_{xx}(a,b) + 2(x-a)(y-b)f_{xy}(a,b) + (y-b)^2f_{yy}(a,b)\Big) + \cdots$$

2nd order coefficients

noniso

m_vis

m_vis

m_vis

m_vis

met

dijetpt

m_vis

m_sv

m_vis

m_vis

ptvis

pt_tt

m_vis

m_vis

m_sv

m_sv

met

m_vis

pt_2

m_sv

m_vis

met

m_vis

pt_tt

met

m_sv

met

ptvis

dijetpt

pt_1

met

met

m_vis

met

pt_tt

dijetpt

mjj

jpt_1

dijetpt

ptvis

0.08

0.08

0.07

jpt_1

pt_2

pt_1

ptvis

pt_tt

met

1.01

0.53

0.50

0.40

0.36

0.27

0.23

0.23

0.21

0.20

0.19

0.18

0.18

0.16

0.15

0.14

0.13

0.13

0.11

0.11

0.11

0.10

0.10

0.10

0.10

0.10

0.09

0.08

0.08

0.08

0.08

0.07

0.07

0.07

0.07

0.06

0.06

0.06

0.06

0.06

0.06

0.06

0.06

0.06

0.06

0.05

0.05

0.05

0.05

0.05

ztt

Application to $H \rightarrow \tau \tau$: 1st order coefficients

ptvis

jpt_1

pt_tt

dijetpt

mjj

met

0.09

0.09

0.09

0.05	0.04	0.07	0.07	0.02	0.10	0 14	0.03	0.10	0.11	0.07	0.04	0 1 9	0.04	017	m_vis m_vis	2.52 1.43	m_vis	m_vis
0.05	0.04	0.07	0.07	0.02	0.10	0.14	0.05	0.10	0.11	0.07	0.04	0.19	0.04	0.17	m_sv ptvis	1.36 0.96	m_sv pt_1	
															pt_tt	0.90	pt_2	
0.15	0.16	0.06	0.08	0.03	0.05	0.23	0.03	0.09	0.12	0.04	0.03	0.52	0.06	0.10	m_vis	0.82	m_vis	
															m_vis	0.79	m_vis	m sv
															IV IV	0.47	jpt_1	ACO V
012	0 1 0	0.00	0.05	0.02	0.05	0.46	0.03	0.20	0.19	0.10	0.03	0.67	0.06	0.00	pt_tt	0.34	m_sv	
0.15	0.10	0.09	0.05	0.02	0.05	0.40	0.05	0.20	0.10	0.10	0.05	0.07	0.00	0.09	m_sv	0.30	ptvis	4
															dijetpt	0.29	pt_tt	nt 1
															m_vis	0.29	m sv	pt_1
0.06	0.04	0.09	0.11	0.02	0.03	0.30	0.02	0.09	0.09	0.22	0.05	0.34	0.08	0.05	dijetpt	0.26		pt_2
															met	0.26	m_sv	
															m_vis	0.25	jpt_2	
0.00	0 1 0	0 1 2	0.07	0.04	014	0.26	0.05	014	0.22	0.24	0.07	0.56	0.09	014	met	0.22	nbtag	Ţ
0.00	0.10	0.15	0.07	0.04	0.14	0.50	0.05	0.14	0.25	0.24	0.07	0.50	0.00	0.14	m_sv	0.21	pt_1	
															ptvis	0.19		
, i-,	N.	÷.	Ň,	, in ,	ъ Б	>	н.	Ś	÷	=	ġ	Ś	ř	ŗ	is	0.19		pt_tt
'ب	ابد	'ب	, L	ايد	ta	N _I	اب	Ξ	Ξı	E	et	.≥'	tp	ne	tt	0.17		met
0	0	ij	ġ	do	q	E	E	ā	o.		þ	's	ije	<u> </u>	pt_tt	0.16	mjj	
				_	-							-	р		m_vis	0.16	pt_1	ntvis
															m_vis	0.15	pt_2	pevib
															mjj	0.14	bpt_1	
• we d	can lea	arn wr	nich va	ariable	es we	re imp	ortan	t for t	ne NN	to ide	entity	tne			m_vis	0.13		jpt_2
resr	pective	e proc	ess												m_sv	0.13	met	
		- p				P/K-t		W-11					VAVE		ptvis	0.12	ptvis	1
					jpt_1	1 1	ptvis	0.11	di	jetpt	dijet	pt	0.09	nbtag	m_vis	0.12	pt_1	
					bpt_1	1	m_vis	0.11	p	tvis	dijet	pt	0.09	nbtag	m_sv	0.12		dijetpt
					nbtag		m_sv	0.11	1	mjj	dijet	pt	0.09	jpt_1	pt_tt	0.11	jpt_2	
					nt ++	met	mii	0.11	יכ	deta nt 2	m_vı	S 2	0.09	ptvis	met	0.11	int 2	Jbt-T
					ipt_1		ipt_2	0.10	J	pt_2 pt	_tt	2	0.09	pt_tt	mii	0.10	jpt_2	
					ptvis	1	otvis	0.10		pt	vis		0.08	jdeta	m_vis	0.10	m_sv	> \ d
					pt_1		pt_2	0.10	I	n_sv	mt_1	L	0.08	ptvis	dijetpt	0.09	pt_1	\sim
					pt_tt	1	pt_tt	0.10	b	pt_1	m_vi	s	0.08		pt_2	0.09	pt_2	$\langle \rangle$
					mjj	nt 2	mjj	0.10	P	t_tt	pt_t	t	0.08	bpt_1	m_sv	0.09	nbtag	F /
					ipt 1	pc_2	pt_tt	0.10	j i	pt_2 pt_1	mjj mii		0.08	bpt.1	m_vis	0.09	ipt 1	
					jpt_1	1	jpt_1	0.10	q	t_tt	dijet	pt	0.08	pt_tt	dijetpt	0.09	pt_1	
					22	jpt_1		0.09		jp	t_2	-	0.08	m_sv	jdeta	0.08	pt_1	d
					ptvis		mii	0.09	т	ot 1	dijet	nt	0.07		dijetpt	0.08	m sv	

50 highest ranked variables or combinations

jpt_1

jpt_1

jpt_1

ptvis

jpt_1

0.07

0.07

0.07

pt_1

pt_1

misc -

noniso

ztt

qqh

ggh

2016, $\tau_h \tau_h$ channel

m_vis

met

m_sv

pt_tt

nbtag

ptvis

jpt_2

jpt_1

mjj

pt_1

m_sv

m_vis

m_sv

pt_tt

ptvis

pt_1

pt_2

jpt_1

m_sv

m_sv

m_vis

m_vis

jpt_1

mjj

jpt_2

m_sv

pt_1

m_sv

nbtag

pt_2

jpt_2

ptvis

m_sv

mt_1

met

jdeta

bpt_1

jpt_1

pt_tt

jpt_1

jpt_2

nbtag

nbtag

jpt_1

ptvis

jpt_1

jpt_1

pt_1

pt_1

pt_1

misc

m_vis

m_vis

m_sv

m_vis

m_vis

m_vis

m_vis

m_vis

ptvis

pt_tt

met

dijetpt

m_sv

m_vis

m_vis

met

m_sv

dijetpt

m_vis

m_sv

m_sv

pt_tt

mjj

m_vis

met

m_vis

m_vis

ptvis

met

met

met

m_sv

met

dijetpt

met

jpt_2

jpt_1

pt_2

met

pt_tt

0.60

0.60

0.25

0.22

0.21

0.21

0.20

0.17

0.17

0.17

0.16

0.16

0.15

0.15

0.14

0.11

0.11

0.11

0.10

0.09

0.09

0.09

0.08

0.08

0.08

0.08

0.07

0.07

0.07

0.07

0.06

0.06

0.06

0.06

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.05

0.04

0.04

0.04

Machine Learning Research for HEP

Some results ... | D. Krücker

Learning for Discovery - The Asimov Loss

Setting the stage

How to optimize for a search?

- We are looking for the best search bin
- Background by some sideband measurement with systematic uncertainty
 - s and b with systematic uncert.
 σ_b= e_{sys} b
- The best statistical answer to this problem: Asymptotic formulae for likelihood-based tests of new physics by CCGV arXiv:1007.1727 [physics.data-an]
 - 1 bin Poisson "expected" discovery significance we call this: "Asimov" significance

Asimov expected discovery significance

Asimov significance looks lengthy

$$Z_A = \left[2\left((s+b) \ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2 + (s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2} \ln\left[1 + \frac{\sigma_b^2 s}{b(b+\sigma_b^2)}\right] \right) \right]^{1/2}$$

- but simplifies without σ_b to $Z_A(\sigma_b = 0) = \sqrt{2((s+b)\ln(1+s/b) s)}$
- and in the limit of small, i.e. s, $\sigma_b \ll b$, to the more common

 $s/\sqrt{b+\sigma_b^2}$

- Here we go for small b and therefore use the full expression above which performs well for even a few events
- NB: Z_A does not scale with Luminosity L

Cross-entropy and Asimov significance

- Neural networks for binary classification are typically trained with cross-entropy
 - Cross-entropy requires the network to model the true label distribution
 - It optimizes accuracy: $Acc = \frac{s_{true} + bt_{rue}}{all}$
- In searches we are not interested in high accuracy, i.e. correct labelling of all background and signal event but
 - We are looking for a phase space region with high signal purity for optimal discovery significance

Cross-entropy and Asimov significance

- When optimizing a classifier a typical approach is to optimize the accuracy: Acc = $\frac{s_{true} + btrue}{all}$
- For neural networks standard approach for training a binary classifier is the **cross-entropy**
- Accuracy maximizing is equivalent to minimizing the cross-entropy:

$$-\frac{1}{N}\sum_{i=1}^{N} \left[y_i^{true} \log(y_i^{pred}) + (1 - y_i^{true}) \log(1 - y_i^{pred}) \right]$$

• Can we optimize directly for the Asimov significance, i.e. can we use it as a loss function ?

Cross-entropy and Asimov significance

- When optimizing a classifier a typical approach is to optimize the accuracy: Acc = $\frac{s_{true} + btrue}{all}$
- For neural networks standard approach for training a binary classifier is the **cross-entropy**
- Accuracy maximizing is equivalent to minimizing the cross-entropy:

• Can we optimize directly for the Asimov significance, i.e. can we use it as a loss function ?

$$Z_A = \left[2\left((s+b) \ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2 + (s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2} \ln\left[1 + \frac{\sigma_b^2 s}{b(b+\sigma_b^2)}\right] \right) \right]^{1/2}$$

The sum is now inside

Cross-entropy and Asimov significance

- When optimizing a classifier a typical approach is to optimize the accuracy: Acc = $\frac{s_{true}+btrue}{all}$
- For neural networks standard approach for training a binary classifier is the **cross-entropy**
- Accuracy maximizing is equivalent to minimizing the cross-entropy:

- Can we optimize directly for the Asimov significance, i.e. can we use it as a <u>loss function</u>?
- **Caveat**: To define the number of signal and background events we need to cut on the discriminator output
 - Makes it non-differentiable ??
 - Differentiability is needed for gradient descent learning

$$Z_A = \left[2\left((s+b) \ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2} \ln\left[1 + \frac{\sigma_b^2 s}{b(b+\sigma_b^2)}\right] \right) \right]^{1/2}$$

Learning for Discovery - Main Idea

Asimov significance as loss function

- A single sigmoid output neuron
- Replace the discrete number of signal and background events by a smooth function of the predicted label y_i^{pred} ∈ [0,1]

$$\begin{split} s &= W_s \sum_{i}^{N_{batch}} y_i^{pred} \times y_i^{true} \\ b &= W_b \sum_{i}^{N_{batch}} y_i^{pred} \times (1 - y_i^{true}) \\ \mathbf{1/Z}_A(s, b) \text{ becomes a smooth} \\ \text{function of } y_i^{pred} \end{split}$$

 W_s, W_b some weights to get the physical number of events depending on luminosity and efficiency to get the

6

Learning for Discovery – applied to a toy SUSY search

Asimov Loss function and Classifier Output - arXiv:1806.00322 [hep-ex]

Asimov loss training

- best Z_A = 6.2 ± 0.6
- Acc = 59%
- AUC=0.80
- Tries to find a background free region

Cross-entropy training + purity cut

- best Z_A = 4.8 ± 0.3
- Acc = 92%
- AUC=0.87

Systematic uncertainty 50%, Differences in Z_A shrink for small systematic uncert.

Asimov score vs. cut on classifier

compressed model point

DESY. Some results ... | D. Krücker

Computing

Some results ... | D. Krücker

CMS Workflow Failures Recovery Panel Towards Al-assisted Operation

Christian Contreras et al. CHEP 2018

Motivation

- CMS central production system
 - ~200K grid cores, ~hundred sites
 - thousands of work flows with thousands of jobs
- Failures are unavoidable
 - Manual operator intervention
 - Can we have an automated system that give advice how to recover a job?
- Train a Deep Neural Network

Trainings data (15,000 workflow tasks)

- Logs (json files) pulled from Workflow Team Web Tools
 which extracts information from Site-Readiness report
 - Job failure codes and site information
 - Operator actions —

- ACDC (rerun failed)
 - splitted
- Clone (total retry)
 - splitted

CMS Workflow Failures Recovery Panel Towards AI-assisted Operation

Christian Contreras et al. CHEP 2018

Difficulties to handle

- Input preprocessing •
 - Getting the data, how to handle missing data
- Unbalanced data •
 - E.g. some failure codes are rare •
 - **SMOTE** Synthetic Minority Over-Sampling by interpolation for minority classes
- Multiclass Training •
 - 2 stage classifier •
 - Binary + multi-class ٠

Model hyper-parameter tuning by Gaussian Processes (skopt)

Eq. number of layers, numer of neurons,

CMS Workflow Failures Recovery Panel Towards Al-assisted Operation

Christian Contreras et al. CHEP 2018

 Confusion matrix
 for the multi-class stage 70/30% split for train/test

1000

750

500

250

0

 A first pass for the supervised learning in error handling prediction. The operator's procedure will be automatized further by applying the decisions that are predicted with acceptable confidence.

Improve current WTC web interface

- To start using Machine Learning Model
 - Include the prediction for recommended action Start recovery from trivial cases
 - Monitor performance for model re-training
- Add GUI display for diagnostic summary reports

Model type	accuracy	Recall	Precision
Binary	87%	87%	83%
Multi-class	87%	86%	83%

DESY. Some results ... | D. Krücker

Collecting data routinely

Ongoing work to implement the tools in the daily workflow for further development

- Operator decisions are "Noisy labels"
- Collecting available data
 - Databases
 - Log files
- Model prediction as additional information to support the operator

AlErrorHandling

- A machine is dedicated to serve as the "action predictor"
- Useful to keep track of several "trained models" during development phase
- Using workflow/task name, extra information are gathered and the "prediction" is produced

Education

Introduction | D. Krücker

Education

DESY provides education in the field of Machine Learning

Schools on Machine Learning within the Terascale Alliance and beyond

- 5th Machine Learning in High Energy ٠ Physics in cooperation with the Yandex Data Science School, June 2019
- 10 days of lab courses in Wuhan/China, ٠ INFIERI summer school, May 2019

1st Terascale School on Machine Learning, • October 2018

2a Train Test

DESY | Intro NN | Dirk Krücker

Many CMS ML activities in different areas