Advances in Machine Learning in experimental High Energy Physics



David Rousseau LAL-Orsay rousseau@lal.in2p3.fr

Uppsala seminar 25th Oct 2017

Outline

1/H

□ ML basics

□ ML in analysis

ML in reconstruction/simulation

- ML challenges
- UWrapping up

Focus on applications rather than details of the techniques

ML in HEP

- Use of Machine Learning (a.k.a Multi Variate Analysis as we call it) already at LEP somewhat, much more at Tevatron (Trees)
- At LHC, Machine Learning used almost since first data taking (2010) for reconstruction and analysis
- □ In most cases, Boosted Decision Tree with Root-TMVA, on ~10 variables
- □ For example, impact on Higgs boson sensitivity at LHC:

analysis	data	no ML	ML	ML
	taking year	sensitivity	sensitivity	data gain
ATLAS $H \to \gamma \gamma$ [16]	2011-2012	4.3	-	-
CMS $H \rightarrow \gamma \gamma \ [17]$	2011-2012	?	2.7	?
ATLAS $H \rightarrow \tau^+ \tau^-$ [18]	2012	2.5	3.4	85%
CMS $H \rightarrow \tau^+ \tau^-$ [19]	2012	3.7	-	-
ATLAS VH \rightarrow bb [20]	2012	1.9	2.5	73%
ATLAS VH \rightarrow bb [21]	2015-2016	2.8	3.0	15%
$\rm CMS \ VH \rightarrow bb \ [22]$	2012	1.4	2.1	125%
$CMS \text{ VH} \rightarrow bb \text{ [23]}$	2015-2016	-	2.8	-

→~50% gain on LHC running

ML in HEP

□ Meanwhile, in the outside world :



- □ "Artificial Intelligence" not a dirty word anymore!
- □ We've realised we're been left behind! Trying to catch up now...

Multitude of HEP-ML events

HiggsML Challenge, summer 2014 → HEP ML NIPS satellite workshop, December 2014 Connecting The Dots, Berkeley, January 2015 Flavour of Physics Challenge, summer 2015 → HEP ML NIPS satellite workshop, December 2015 □ DS@LHC workshop, 9-13 November 2015 LHC Interexperiment Machine Learning group Started informally September 2015, gaining speed IML workshop @CERN 20-22 March 2017 Moscou/Dubna ML workshop 7-9th Dec 2015 Heavy Flavour Data Mining workshop, 18-21 Feb 2016 Connecting The Dots, Vienna, 22-24 February 2016 Hep Software Foundation workshop 2-4 May 2016 at Orsay, N Connecting The Dots, LAL-Orsay, 6-9 March 2017 DS@HEP workshop @FNAL 8-12 May 2017 ACAT conference Seattle, Sep 2017 Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oc

November 2015, CER **C**ennecting Bezio the **Dots** unversité Intelligent Trackers 2017 6th - 9th March 2017, AL-Orsay, France

Goh Energy Physics

ML Basics



BDT in a nutshell



- □ Single tree (CART) <1980
- □ AdaBoost 1997 : rerun increasing the weight of misclassified entries → Boosted Decision Trees (Gradient BDT, random forest...)

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Classifier basics



Neural Net in a nutshell



- Neural Net ~1950!
- But many many new tricks for learning, in particular if many layers (also ReLU instead of sigmoïd activation)
- "Deep Neural Net" up to 100 layers
- Computing power (DNN training can take days even on GPU) Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Deep learning



No miracle

- ML (nor Artificial Intelligence) does not do any miracles
- For selecting Signal vs Background and underlying distributions are known, nothing beats Likelihood ratio! (often called "bayesian limit"):
 - $O L_{S}(x)/L_{B}(x)$
- OK but quite often L_S L_B are unknown

+ x is n-dimensional

- ML starts to be interesting when there is no proper formalism of the pdf
- ➡ mixed approach, if you know something, tell your classifier instead of letting it guess

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017



ML Techniques



Overtraining



Evaluated on independent test dataset (correct)

Score distribution different on test dataset wrt training dataset → "Overtraining"== possibly excessive use of statistical fluctuation



Standard basic way (default TMVA until recently)

under/over training



Two-fold Cross Validation

IIII



- \rightarrow test statistics = total statistics
- →double test statistics wrt one fold CV
- \rightarrow (double training time of course)

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

5-fold Cross Validation



same test statistics wrt two-fold CV, larger training statistics 4/5 over ½ (larger training time as well)

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

101 75 a 14

A

5-fold Cross Validation

MITT



71.70

14

A

5-fold Cross Validation

IIIII



101

14

A

5-fold Cross Validation

IIIII



5-fold Cross Validation



Note : if hyper-parameter tuning, need a third level of independent sample "nested CV"

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

5-fold Cross Validation "à la Gabor"



Average of the scores on A B C D is **often** better than the score of one training ABCD bonus: variance of the samples an estimate of the statistical uncertainty (also save on training time) Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017



71-1

- →Testing variance
- →Testing variance
- →Testing variance
- →Testing variance
- → Testing variance

What does a classifier do?



The classifier "projects" the two multidimensional "blobs" maximising the difference, without (ideally) any loss of information

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Re-weighting

Suppose a variable distribution is slightly different between a Source (e.g. Monte Carlo) and a Target (e.g. real data)

 \circ \rightarrow reweight! ...then use reweighted events



- □ What if multi-dimension ?
- Usually : reweight separately on 1D projections, at best 2D, because of quick lack of statistics
- Can we do better ?

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017



Multi dimensional reweighting (2)

- Reweighting the Source distribution on the score allows multidimensional reweighting without statistics problem
- Usual caveat still hold : Target support should be included in Source support, distributions should not be too different otherwise unmanageable very large or very small weights
- (Note : "reweighting" in HEP language <==> "importance sampling" in ML language)

Anomaly : point level

- Also called outlier detection
- Two approaches:
 - Unsupervised : give the full data, ask the algorithm to cluster and find the lone entries : o1, o2, O3



X

Supervised : we have a training "normal" data set with N1 and N2.
 Algorithm should then spot o1,o2, O3 as "abnormal" i.e. "unlike N1 and N2" (no a priori model for outliers)

Application : detector malfunction, grid site malfunction, or even new physics discovery... Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Anomaly : population level

- Also called collective anomalies
- Suppose you have two independent samples A and B, supposedly statistically identical. E.g. A and B could be:
 - MC prod 1, MC prod 2
 - MC generator 1, MC generator 2
 - Geant4 Release 20.X.Y, release 20.X.Z
 - Production at CERN, production at BNL
 - Data of yesterday, Data of today
- How to verify that A and B are indeed identical ?
- Standard approach : overlay histograms of many carefully chosen variables, check for differences (e.g. KS test)
- One ML approach (not the only one): ask an artificial scientist, train your favorite classifier to distinguish A from B, histogram the score, check the difference (e.g. AUC or KS test)
 - \rightarrow only one distribution to check



Local big difference (e.g. non overlapping distribution, hole)



Advances in ML in HESCORVED Rousseau, Uppsala seminar, 25 Oct 2017 8B 30

HSF ML RAMP on anomaly

- RAMP : collaborative competition around a dataset and a figure of merit. Organised in June 2016 by CDS Paris Saclay with HEP people. See <u>agenda.</u>
 - Dataset built from the Higgs Machine Learning challenge dataset (on CERN Open Data Portal)
 - Lepton, and tau hadron 3 momentum, MET : PRImary variables
 - DERived variables e.g various invariant masses (computed from the above) from Htautau analysis
 - o →reference dataset
 - Skewed" dataset built from the above, introducing small and big distortions:
 - Change of tau energy scale (Small scaling of Ptau)
 - Holes in eta phi efficiency map of lepton and tau hadron
 - Outliers introduced, each with 5% probability
 - Eta tau set to large non possible values
 - P lepton scaled by factor 10
 - Missing ET + 50 GeV
 - Phi tau and phi lepton swapped → DERived variables inconsistent with PRImary one
 - o → skewed dataset

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

HSF ML RAMP on anomaly (2)



Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

ML Tools



Modern Software and Tools

- New version of TMVA (root 6.0.8 on beyond) (see talk Lorenzo Moneta, Sergei Gleyzer IML workshop CERN March 2017)
 - o Jupyter interface
 - Hyper-parameter optimisation
 - o Cross-validation
 - o (...unfortunately not so well documented yet)
- Non HEP software
 - Sci-kit learn : de facto standard toolbox ML (except Deep Learning) (python, but fast)
 - Keras+Thenao/TensorFlow : NN toolbox (build a NN in a few lines of python)
 - XGBoost best BDT on the market, both speed and performance (c++ with python interface)
- Note : for ~10 variable classification/regression task gradient BDT is still the tool of choice!
- Platforms
 - Your laptop is sufficient in many cases : install e.g. Anaconda <u>https://docs.continuum.io/anaconda/install</u> (<u>demo</u>)
 - If not, more and more platforms looking for users, maybe on your campus (with GPU DNN ==millions of parameter to optimise=>heavy duty linear algebra)
 - o 50 GPU platform at Lyon CC-IN2P3, little used so far
- For CERN users:
 - o SWAN interactive data analysis on the web see <u>https://swan.web.cern.ch/content/machine-learning</u>
 - CVMFS ML setup for any CVMFS enabled platform

ML in analysis



Candidat H→Z(→μ⁺μ⁻)Z(→e⁺e⁻)

Run Number: 182796, Event Number: 74566644 Date: 2011-05-30, 06:54:29 CET

EXPERIMEN

EtCut>0.3 GeV PtCut>2.0 GeV Vertex Cuts: Z direction <1cm Rphi <1cm

Muon: blue Electron: Black Cells: Tiles, EMC

Deep learning for analysis



Signal efficiency

Advances in ML in HEP, David Rousseau, Up

Deep learning for analysis (2)

1410.3469 Baldi Sadowski Whiteson

□ H tautau analysis at LHC: H→tautau vs Z→tautau

- Low level variables (4-momenta)
- High level variables (transverse mass, delta R, centrality, jet variables, etc...)



- Here, the DNN improved on NN but still needed high level features
- Both analyses with Delphes fast simulation
- ~10M events used for training (>>10* full G4 simulation in ATLAS)

Systematics-aware training

Our experimental measurement papers typically ends with

• measurement = m $\pm \sigma$ (stat) $\pm \sigma$ (syst)

- o σ (syst) systematic uncertainty : known unknowns, unknown unknowns...
- □ Name of the game is to minimize quadratic sum of :

 σ (stat) $\pm \sigma$ (syst)

- \Box ML techniques used so far to minimise σ (stat)
- □ Impact of ML on σ (syst) or even better global optimisation of σ (stat) ± σ (syst) is an open problem
- \Box Worrying about σ (syst) untypical of ML in industry
- However, a hot topic in ML in industry: transfer learning
- E.g. : train image labelling on a image dataset, apply on new images (different luminosity, focus, angle etc...)
- □ For HEP : we train with Signal and Background which are not the real one (MC, control regions, etc...)→ source of systematics^{Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017}

Syst Aware Training: adversarial

3

Inspired from 1505.07818 Ganin et al :

ACAT 2017 Ryzhikov and Ustyuzhanin





Parameterised learning

<u>1601.07913</u> Baldi, Cranmer, Faucett, Sadowksi, Whiteson





Parameterised learning (2)



- Train on 28 features plus true mass
- Parameterised NN as good as single mass training
- ❑ → clean interpolation
- (mass just an example)
- Very recently used by CMS bblvl v search <u>https://arxiv.org/pdf/1708.0</u> <u>4188.pdf</u>

ML in reconstruction



Jet Images

arXiv 1511.05190 de Oliveira, Kagan, Mackey, Nachman, Schwartzman

- Distinguish boosted W jets from QCD
- Particle level simulation
- Average images:











240 < p,/GeV < 260 GeV, 65 < mass/GeV < 95 Pythia 8, QCD dijets, $\sqrt{s} = 13 \text{ TeV}$ [GeV] € [Translated] Azimuthal Angle Pixel p 10 0.5 0.5 10-1 10⁻² 10⁻³ 10-4 10⁻⁵ -0.5 10⁻⁶ 107 10⁻⁸ -0.5 0 0.5 -1 [Translated] Pseudorapidity (n)



Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

45

Jet Images : Convolution NN



RNN for b tagging

- BDT and usual NN expect a fix number of input. What to do when the number of inputs is not fixed like the tracks for b-quark jet tagging ?
- Recurrent neural networks have seen outstanding performance for processing sequence data
 - Take data at several "time-steps", and use previous time-step information in processing next time-steps data

ATL-PHYS-PUB-2017-003

- For b-tagging, take list of tracks in jet and feed into RNN
 - Basic track information like d0, z0, pt-Fraction of jet, ...
 - Physics inspired ordering by d0-significance
- RNN outperforms other IP algorithms
 - No explicit vertexing, still excellent performance
 - First combinations with other algorithms in progress
- Learning on sequence data may be important in other places!





ML in simulation



Generative Adversarial Network



Condition GAN

Text to image

MITT

this small bird has a pink breast and crown, and black primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma

this magnificent fellow is almost all black with a red crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen





GAN for simulation



CaloGAN



Paganini et al.





- \Box σ_1 :width in Middle layer
- One of many physics variable examined
- Pion more difficult
- J →very promising

Data Challenges



Higgs Machine learning challenge

See talk DR CTD2015 Berkeley

- An ATLAS Higgs signal vs background classification problem, optimising statistical significance
- Ran in summer 2014
- 2000 participants (largest on Kaggle at that time)
- Outcome
 - Best significance 20% than with Root-TMVA
 - (gradient) BDT algorithm of choice in this case where number variables and number of training events limited (NN very slightly better but much more difficult to tune)
 - XGBoost written for HiggsML, now best BDT on the marke
 - Wealth of ideas, documented in <u>JMLR proceedings v42</u>
 - Still working on what works in real life what does not
 - Raised awareness about ML in HEP

Also:

- Winner Gabor Melis hired by DeepMind
- Tong He, co-developper of XGBoost, winner of special "HEP meets ML" price got a PhD grant and US visa



Towards a Future Tracking Machine Learning challenge



A collaboration between ATLAS and CMS physicists, and Machine Learners



TrackML : Motivation

- □ See details DR talk at CTD/WIT 2017
- Tracking (in particular pattern recognition) dominates reconstruction CPU time at LHC
- □ HL-LHC (phase 2) perspective : increased pileup :Run 1 (2012): <>~20, Run 2 (2015): <>~30, Phase 2 (2025): <>~150
- CPU time quadratic/exponential extrapolation (difficult to quote any number)
- Large effort within HEP to optimise software and tackle micro and macro parallelism. Sufficient gains for Run 2 but still a long way for HL-LHC.
- >20 years of LHC tracking development. Everything has been tried?
 - Maybe yes, but maybe algorithm slower at low lumi but with a better scaling have been dismissed ?
 - Maybe no, brand new ideas from ML (i.e. Convolutional NN)





Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017





TrackML : engaging Machine Learners

- Suppose we want to improve the tracking of our experiment
- We read the literature, go to workshops, hear/read about an interesting technique (e.g. ConvNets, MCTS...). Then:
 - Try to figure by ourself what can work, and start coding traditional way
 - Find an expert of the new technique, have regular coffee/beer, get confirmation that the new technique might work, and get implementation tips **>** better
- ...repeat with each technique...

Much much better:

- Release a data set, with a benchmark, and have the expert do the coding him/herself
- → he has the software and the know-how so he'll be (much) faster even if he does not know anything about our domain at the beginning
- →engage multiple techniques and experts simultaneously (e.g. 2000 people participated to the Higgs Machine Learning challenge) in a comparable way
- o →even better if people can collaborate
- \rightarrow a challenge is a dataset with a benchmark and a buzz
- Looking for long lasting collaborations beyond the challenge
- Focus on the pattern recognition : release list of 3D points, challenge is to associate them into tracks fast. Use public release of ATLAS tracking (<u>ACTS</u>) as a simulation engine and starting kit
- Phase 1 (just accuracy) will run winter 2018
- □ Phase 2 (accuracy and CPU) will run summer 2018

Pattern recognition



Real-time face recognition : efficiency, fake, CPU time...

Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Pattern Recognition/Tracking

- Pattern recognition/tracking is a very old, very hot topic in Artificial Intelligence, but very varied
- Note that these are real-time applications, with CPU constraints







in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

A recent attempt : NOVA



CTDWIT 2017 2D tracking Hackathon

- CTDWIT 6-9th March 2017 LAL-Orsay
- Very simplified 2D simulation with HL-LHC ATLAS layout (circular detectors, multiple scattering, inefficiency, stopping tracks)
 EPJ Web Conf., 150 (2017) 00015
- Run on RAMP platform
- 30 people (tracking experts mostly) for 2 hours in the same room, plus 36 hours till the end of the conference
- □ Winner is a Monte Carlo Tree Search algorithm (used in Go algorithms before and also by Alpha-Go)
- Runner-up a "real" ML algorithm : Long Short Term Memory



Belle II Experiment @belle2collab · 15 min

Congrats to four **#Belle2** PhD students for winning the Tracking Challenge at this year's Connecting the DotsD Conference! **#ctdwit #hackathon**

À l'origine en anglais



David Rousseau

.@SteveAFarrell winner of #CTDWIT TrackMLRamp 2D #hackathon at @LALOrsay in the ML category. Congrats !

À l'origine en anglais



Wrapping-up



More on ML in HEP history

Computer Physics Communications 49 (1988) 429-448 North-Holland, Amsterdam

NEURAL NETWORKS AND CELLULAR AUTOMATA IN EXPERIMENTAL HIGH ENERGY PHYSICS

B. DENBY

Laboratoire de l'Accélérateur Linéaire, Orsay, France

Received 20 September 1987; in revised form 28 December 1987

- 1987 Very first ML in HEP paper known
- ML for tracking and calo clustering
- B. Denby then moved from Delphi at LEP to CDF at Tevatron. He still active outside HEP: 2017 analysis of ultrasonic image of the tongue
- 1992 JetNet Carsten Peterson, Thorsteinn Rognvaldsson (Lund U.), Leif Lonnblad (CERN) (~500 citations) really started NN use in HEP

Advances in ML in HEP, David Rousseau, Uppsala ser



ML playground

IIII



Advances in ML in HEP, David Rousseau, Uppsala seminar, 25 Oct 2017

Collection of links

- In addition to workshops mentioned in the first transparencies, and references mentioned in the talks
- Interexperiment Machine Learning group (IML) is gathering speed (documentation, tutorials, etc...). Topical monthly meeting. Workshop 20-22 March :
- □ An internal ATLAS ML group has started in June 2016. In CMS in June 2017
- https://higgsml.lal.in2p3.fr
- <u>http://opendata.cern.ch/collection/ATLAS-Higgs-Challenge-2014</u>: permanent home of the challenge dataset
- NIPS 2014 workshop agenda and proceedings http://jmlr.org/proceedings/papers/v42/
- Mailing list opened to any one with an interest in both Data Science and High Energy Physics : <u>HEP-data-science@googlegroups.com</u> and <u>Ihc-machinelearning-</u> wg@cern.ch
- IN2P3 project starting <u>http://listserv.in2p3.fr/cgi-bin/wa?A0=MACHINE-LEARNING-L</u> open to anyone with some interest to ML (planning on 2 x 1day workshop per year)
- □ IN2P3 School of Statistics 28 May 1 June 2018 To be Confirmed (see <u>SoS 2016</u>)

ML Collaborations

- Many of the new ML techniques are complex→difficult for HEP physicists alone
- □ ML scientists (often) eager to collaborate with HEP physicists
 - o prestige
 - o new and interesting problems (which they can publish in ML proceedings)
- Takes time to learn common language
- Access to experiment internal data an issue, but there are ways out
- Note : Yandex Data School of Analysis (with ~10 ML scientists) now a bona fide institute of LHCb
- Very useful/essential to build HEP ML collaborations : study on shared dataset, thesis (Computer Science or HEP)
- □ There is always a friendly Machine Learner on a campus!

Open Data

- Public dataset are essential to collaborate (beyond talking over beer/coffee) on new ML techniques with ML experts (or even physicists in other experiments)
 - o can share without experiments Non Disclosure policies
- Some collaborations built on just generator data (e.g. Pythia) or with simple detector simulation e.g. Delphes
 - o good for a start, but inaccurate
- Effort to have better open simulation engine (e.g. Delphes 4-vector detector simulation, ACTS for tracking)
- □ UCI dataset repository has some HEP datasets
- Role of CERN Open Data portal:
 - We (ATLAS) initially saw its use for outreach purposes (CMS has been more open on releasing data)
 - But after all, ML collaboration is a kind of scientific outreach
 - →ATLAS uploaded there in 2015 the data from Higgs Machine Learning challenge (essentially 4-vectors from full G4 ATLAS simulation Higgs->tautau analysis)
 - ATLAS consider releasing more datasets dedicated to ML studies

Conclusion

- We (in HEP) are analysing data from multi-billion € projects→should make the most out of it!
- Recent explosion of novel (for HEP) ML techniques, novel applications for Analysis, Reconstruction, Simulation, Trigger, and Computing
- Some of these are ~easy, most are complex: open source software tools are ~easy to get, but still need (people) training, know-how
- More and more open datasets/simulators
- More and more HEP and ML workshops, forums, schools, challenges
- More and more direct collaboration between HEP researchers and ML researchers
- □ HEP will need more and more access to (GPU) training resources
- Never underestimate the time for :
 - o (1) Great ML idea→
 - (2) ...demonstrated on toy dataset →
 - o (3) ... demonstrated on real experiment analysis/dataset \rightarrow
 - o (4) ... experiment publication using the great idea