

Machine Learning



meets

particle physics

Overview & examples

Troels C. Petersen

With contributions from Stefan Hasselgren, Lukas Ehrke, Frederik Faye, Benjamin Henckel, and Daniel Nielsen Niels Bohr Institute

AI center seminar, 3rd of May 2019



Machine Learning



meets

particle physics

Overview & examples



Troels C. Petersen

With contributions from Stefan Hasselgren, Lukas Ehrke, Frederik Faye, Benjamin Henckel, and Daniel Nielsen Niels Bohr Institute

AI center seminar, 3rd of May 2019

Outline

Outline of talk:

- In the beginning: **b-jet tagging in ALEPH**
- HEP data and why it is exciting for Machine Learning (ML)
- Going large scale: Electrons and photons in ATLAS
 - Data samples, variables, and selections
 - Electron PID and Energy Reconstruction (ER)
 - Discussion of performance measures (loss functions)
- Looking at the future: v-reconstruction in IceCube

Purpose of talk:

- Show Machine Learning cases in science.
- Present HEP data, and why it is great (but also cumbersome!)
- Open up for possible inspiration / collaboration

In the following, all numbers and plots are "Not Even Preliminary", and should in not be used elsewhere.

A word about particle physics

We search for MANY different things, typically rare (1:10⁹) with complex decays.



Particle physics data and simulation

To make sure that we understand our experiment we use simulation extensively:

- Detector optimisation (before experiment)
- Reconstruction design/optimisation (before/during experiment)
- Selection optimisation (during experiment)
- Signal efficiency estimates (during experiment)



The simulation is done in three steps: Event generation, Simulation, and Digitisation.

The total CPU time needed for one event is about 20-30 minutes, and we have now simulated about a billion events (using 0.5M cores).

The simulation is done from **first principles**, and there are therefore (smaller) **differences between simulation and data** (maybe a point of interest to fix?).

Aim of this project

Electrons and photons play a central role in the ATLAS physics programme, in particular in **Higgs physics**, where they dominate the two golden channels. *Current methods use likelihood approach (PID) and simple ML (Energy)*.

Given the cost of running, we would be satisfied, if we could add 10% to each of these in terms of statistics, knowing this would also benefit many other analysis.





25 years ago, particle physics was actually at the forefront of using Machine Learning. We had large computers and much data fit for ML usage.

At the time, LEP was searching for the Higgs boson at lower masses, where its decay was almost always to b-quarks.

For this reason, many resources were used to get the best possible b-jet tagging in place.

25 years ago, particle physics was actually at the forefront of using Machine Learning. We had large computers and much data fit for ML usage.

At the time, LEP was searching for the Higgs boson at lower masses, where its decay was almost always to b-quarks.

For this reason, many resources were used to get the best possible b-jet tagging in place.

Both lifetime (displaced vertices), jet shape, and lepton pT was used, but non of these by themselves provide a good way to select b-jets.



Impact parameter significance:

Figure 3.5: Impact parameter significance (δ/σ_{δ}) distributions for (a) fragmentation or tracks from uds events, (b) tracks from weakly decaying c hadrons, (c) tracks directly from b hadron decay, (d) tracks from the cascade decay of a c hadron from the decay of a b hadron.

However, using one of the very first ML algorithms (JetNet 3.4), six variables were put together in a neural network with two hidden layers each with 10 neurons:

- Light quark (uds) jet probability from track impact parameter significance.
- Difference in Chi2 from search for secondary vertices in jet.
- Transverse momentum of (possible) electron/muon in jet.
- Boosted sphericity of jet.
- Energy flow multiplicity (scaled by jet energy).
- Sum of transverse momenta (with respect to the jet axis) squared.

The neural network was trained on 400.000 **simulated** events, and though I haven't been able to find the exact time used for this training, colleagues have told me "many hours, sometime days".

Interestingly, my students now code the setup in about an hour, and get results in minutes.



The result of these labours was a very nice b-jet tagging variable, which allowed ALEPH to get the most out of their data.







Overview

Currently in ATLAS, electrons and photons are identified using a likelihood approach:

 $4 \,\mathrm{ma}$

C	п				
$d_{\mathcal{L}} = \frac{\mathcal{L}_S}{\mathcal{L}_S + \mathcal{L}_B},$	$\mathcal{L}_{S}(\vec{x}) = \prod P_{s,i}(x_i)$				
	i=1				

The likelihood is composed of 22 variables, for which 1D histograms are used to compute the likelihood value. 9 calor

To minimise correlations, the likelihood is divided into regions of E_T and η .

This makes for a very transparent approach, which at the same time performs well.

The question is, if there is more information to be gained, and thus a more powerful PID to be gotten.

Enter Machine Learning (ML)...

	Туре	Description				
1	Hadronic leakage	Ratio of $E_{\rm T}$ in the first layer of the hadronic calorimeter to $E_{\rm T}$ of the EM cluster				
24		(used over the range $ \eta < 0.8$ or $ \eta > 1.37$)				
1.20		Ratio of $E_{\rm T}$ in the hadronic calorimeter to $E_{\rm T}$ of the EM cluster	R _{Had}			
100		(used over the range $0.8 < \eta < 1.37$)				
1	Back layer of	Ratio of the energy in the back layer to the total energy in the EM accordion	f_3			
	EM calorimeter	calorimeter				
	Middle layer of	Lateral shower width, $\sqrt{(\Sigma E_i \eta_i^2)/(\Sigma E_i) - ((\Sigma E_i \eta_i)/(\Sigma E_i))^2}$, where E_i is the	$W_{\eta 2}$			
	EM calorimeter	energy and η_i is the pseudorapidity of cell <i>i</i> and the sum is calculated within				
-		a window of 3×5 cells				
		Ratio of the energy in 3×3 cells over the energy in 3×7 cells centered at the				
		electron cluster position				
imeter	variables	Ratio of the energy in 3×7 cells over the energy in 7×7 cells centered at the				
	1	electron cluster position				
	Strip layer of	Shower width, $\sqrt{(\Sigma E_i(i - i_{\text{max}})^2)/(\Sigma E_i)}$, where <i>i</i> runs over all strips in a window	wstot			
	EM calorimeter	of $\Delta \eta \times \Delta \phi \approx 0.0625 \times 0.2$, corresponding typically to 20 strips in η , and				
		$i_{\rm max}$ is the index of the highest-energy strip				
1		Ratio of the energy difference between the largest and second largest energy	$E_{\rm ratio}$			
-		deposits in the cluster over the sum of these energies				
1.1		Ratio of the energy in the strip layer to the total energy in the EM accordion	f_1			
13		calorimeter				
1.5	Track quality	Number of hits in the B-layer (discriminates against photon conversions)	n _{Blayer}			
1.19		Number of hits in the pixel detector	n _{Pixel}			
3		Number of total hits in the pixel and SCT detectors	n _{Si}			
cking variables		Transverse impact parameter	d_0			
		Significance of transverse impact parameter defined as the ratio of d_0				
		and its uncertainty				
		Momentum lost by the track between the perigee and the last	$\Delta p/p$			
		measurement point divided by the original momentum				
	TRT	Total number of hits in the TRT				
15		Ratio of the number of high-threshold hits to the total number of hits in the TRT				
1000	Track-cluster	$\Delta \eta$ between the cluster position in the strip layer and the extrapolated track	$\Delta \eta_1$			
	matching	$\Delta \phi$ between the cluster position in the middle layer and the extrapolated track	$\Delta \phi_2$			
1.		Defined as $\Delta \phi_2$, but the track momentum is rescaled to the cluster energy				
ching	variables	before extrapolating the track to the middle layer of the calorimeter				
		Ratio of the cluster energy to the track momentum	E/p			
	Conversions	Veto electron candidates matched to reconstructed photon conversions	isConv			

1 conversion variable (for photons)



Reweighing



Reweighing



Electron PID performance

The electron PID performance is generally much improved with ML:





We train the Machine Learning (ML) algorithm (LightGBM) with a mix of backgrounds, and then see how well it performs on each. We compare to the current ATLAS LH, not to boast our results, but as a solid reference, which helps us getting the most performant & generel results.



Electron PID performance

The electron PID performance is generally much improved with ML:



Where do we improve (most)?

The improvements are NOT uniform in energy and angle. We gain most in the "crack" and forward direction.



Electron PID feature importance

The importance of each PID input variable is show below (https://github.com/slundberg/shap).



Electron PID feature importance

The importance of each PID input variable is show below (https://github.com/slundberg/shap).









Electron PID on data



Impact in data

ML electron PID on probe side yields in data more Zee events (same background):



Data: Label confusion

In real data we don't have perfect labelling. The selections for signal and background have a few percent of label confusion each, depending on energy and detector part hit.

It turns out, that Random Forests (RF) perform better than Boosted Decision Trees (BDT) given this label confusion, as one might expect.

LightGBM RF with 100 trees - AUC value					LightGBM BDT with 100 trees - AUC value				
0.9966	0.9969	0.9968	0.9968	0.9969	0.9928	0.9918	0.9896	0.9853	0.9612
0.9967	0.9968	0.9967	0.9967	0.9969	0.9933	0.9936	0.9924	0.9904	0.9807
0.9967	0.9967	0.9967	0.9967	0.9968	0.9939	0.9934	0.9928	0.9918	0.983
0.9966	0.9966	0.9967	0.9966	0.9967	0.9937	0.9935	0.9928	0.9926	0.9855
0.9967	0.9966	0.9968	0.9967	0.9968	0.9936	0.9936	0.9934	0.9926	0.9862
0	1.0 Mislabelir	5.0 ng of backg	10.0 round [%]	25.0	0	1.0	5.0	10.0	25.0



Electron ER - BDT vs. CNN

We started to work on electron energy reconstruction (ER) using scalar variables combined with a BDT approach, just like ATLAS does.

However, we are now exploring to use a Convoluted Neural Network (CNN) for the task, as this "naturally" fits the problem, when considering the calorimeter cells as images.

Naturally, there are still scalar variables to add to the regression:



CNN scalar variables

- p deltaPhiRescaled2
- pX_deltaPhiFromLastMeasurement
- pX deltaPhiRescaled0
- pX deltaEta2
- pX deltaEta3
- p_charge
- BC distanceFromFront
- BC filledBunches
- p pt track
- averageInteractionsPerCrossing
- NvtxReco
- No ECAL variables

Photon ER performance



Photon ER performance



CNN setup

- ECAL layers 0 (presampler), 1, 2 and 3 (cell values and times)
- Frederik Faye
- All images are produced to be 56×55 (better than 7×11 and 56×11) and are stacked to give a $56 \times 55 \times 4$ array
- Cell values are scaled to GeV
- Time values are divided by 50 ns

NOTE: For now, we only consider barrel electrons ($|\eta| < 1.3$)



The great thing is that for each cell we don't just have the energy, but also the time (rejecting out-of-time pile-up), gain, and cell noise level (gauging the energy precision).

However, these are not same units, so combined with gate (not concatenation).

CNN architecture

We use a 3×3 convolution matrices for all layers.

Each convolution layer is followed by a batch normalisation and activation.

For all i > 1, block begins with downsampling and the number of feature maps is doubled.

A worry is, that the scalar variables "drown" in the many feature map outputs. To be investigated further.

However, we know that scalar variables improve performance as it is!

Images containing time are treated differently...



20

Frederik Faye



J

20

Frederik Faye



5

20





20



Electron Energy Reconstruction

- on MC... latest development!

FiLM = Feature wIse Linear Modulation

 $\operatorname{FiLM}(\mathbf{x}) = \gamma(\mathbf{z}) \odot \mathbf{x} + \beta(\mathbf{z}).$



https://distill.pub/2018/feature-wise-transformations/

Results from FiLM





The IceCube detector is a less "classic" particle physics detector. Here, 86 strings with about 5000 Digital Optical Modules (DOMs) in total are put in the ice at the South Pole, and used to detector neutrinos (and involuntarily cosmic muons) interact in the ice.

The detector is triggered by coincidences of several adjacent DOMs, and then read out.

Each DOM provides a measurement in time and size of signal. However, there is a significant amount of noise and also effects such as after-pulses, which makes the data less clean.



The IceCube detector is a less "classic" particle physics detector. Here, 86 strings with about 5000 Digital Optical Modules (DOMs) in total are put in the ice at the South Pole, and used to detector neutrinos (and involuntarily cosmic muons) interact in the ice.

The detector is triggered by coincidences of several adjacent DOMs, and then read out.

Each DOM provides a measurement in time and size of signal. However, there is a significant amount of noise and also effects such as after-pulses, which makes the data less clean.

The **bottleneck** is the event reconstruction!

This is based on the minimisation of a likelihood including ray tracing and ice properties.



Neutrinos and cosmic muons interact in the ice, and leaves signals to be reconstructed.

Neutrinos and cosmic muons interact in the ice, and leaves signals to be reconstructed.

Problem 1: Which hits belong to the event and which are noise?

Neutrinos and cosmic muons interact in the ice, and leaves signals to be reconstructed.

Problem 1: Which hits belong to the event and which are noise?

Problem 2: Given a list of hits, how to determine the direction, energy, type, etc.? And... how to do it in a "reasonable" amount of time? Currently t(reco) = 30 min.

Neutrinos and cosmic muons interact in the ice, and leaves signals to be reconstructed.

Problem 1: Which hits belong to the event and which are noise?

A student of mine (Andreas Søgaard) tried to see, if he could get an ML algorithm to do the reconstruction. It didn't perform very well (yet!), **but t(reco) = 0.01 sec.** Problem 2: Given a list of hits, how to determine the direction, energy, type, etc.? And... how to do it in a "reasonable" amount of time? Currently t(reco) = 30 min.

Conclusions

I think that there is a lot of prospect in Machine Learning for physics.

- New algorithms see the light of day almost daily.
- In some cases, it may simply give a more performant data analysis.
- However, in some cases, it makes all the difference.
- Particle physics data is well suited for ML as we have "accurate" simulation.

The data requires collaboration, as there are several "particle physics tricks" needed to evaluate performance in real data.

There are many areas to try ML on:

- Transformation of simulation to match data better (challenge: extrapolation).
- Simulation using GANs and VAEs (already started in ATLAS).
- Reconstruction in the IceCube experiment

I've been surprised by the speed with which students "pick up" ML, once you give them an introduction to it. The challenge is often to find data "suitable" for the algorithms given. We - in HEP - tend to actually have that!



Signal/background selection

- DAOD production: mixture of EGAM1, EGAM3, EGAM7, EGAM8 and EGAM9 including cells and bunch crossing information
- $Z \rightarrow ee$ Tag and Probe
 - Tag
 - $p_T > 24.5 \, \text{GeV}$
 - Tight ID
 - Loose Isolation
 - Crack veto and Central
 - Pass HLT_e26_Ihtight_nod0_ivarloose
 - Has track particle and vertex
 - Pass object quality cut

- Probe
 - $p_T > 4.5 \, {\rm GeV}$
 - Pass object quality cut
 - $|\eta| < 4.9$
- Combined
 - dR > 0.4 for tag and probe
 - *M_{ee}* > 50 GeV
- Dijet, $W^{\pm} \rightarrow \ell^{\pm} \nu \ (\ell = e, \mu), Z \rightarrow \mu \mu$ samples background selection
 - Missing transverse energy (MET) $< 25 \, \mathrm{GeV}$
 - $p_T > 4.5 \, \text{GeV}$
 - pass Object Quality
 - Z veto: Match with any other medium electron, $|M_{ee} m_{Z^0}| < 20 \text{ GeV}$
 - W veto: Match with MET, transverse mass < 40 GeV
- $Z(\gamma) \rightarrow ee$ (Drell-Yan), $Z \rightarrow \ell \ell \gamma$, $\gamma + jet$
 - EGamma truth particles
- More complete documentation https://twiki.cern.ch/twiki/bin/ viewauth/AtlasProtected/EgammaMachineLearning

Purities after T&P

The electron probe purities for both signal and background are far from ideal!

Signal sample purities										
	η : 2.01-2.47 ·	58.8%	83.2%	95.4%	98.2%	96.3%	0.0%		η : 2.0	
	η : 1.52-2.01 ·	52.5%	80.1%	95.4%	98.2%	96.1%	90%		η : 1.5	
	η : 1.37-1.52 ·	45.7%	75.8%	94.4%	97.9%	95.5%	- 80%		η : 1.3	
	η : 0.8-1.37 ·	47.9%	78.9%	95.9%	98.4%	95.9%	-70%		η : 0	
	η : 0.0-0.8 ·	47.5%	79.9%	96.0%	98.5%	96.0%	- 60%		n :	
	Er/Gev Er/Gev Er/Gev Er/Gev : 20.30									

Background sample purities									
η : 2.01-2.47 -	70.0%	69.3%	83.4%	89.1%	91.2%	000%			
η : 1.52-2.01 ⁻	75.2%	77.5%	88.2%	92.1%	94.3%	- 90%			
η : 1.37-1.52 ·	66.1%	71.3%	85.2%	89.5%	91.9%	- 80%			
η : 0.8-1.37 ⁻	74.4%	80.8%	91.1%	94.6%	96.7%	- 70%			
η : 0.0-0.8 ⁻	80.6%	83.0%	91.8%	95.0%	97.0%	- 60%			
ETIGev ETIGev ETIGev ETIGev.									

Comparison of Combination



Performance of all methods



Electron PID on data



What is a CNN?

CNNs are a type of neural network, which works well with spatially dependent data (typically images). CNNs use parameter/weight sharing.

Multi-layered images (e.g. RGB or ATLAS calorimeter) are handled naturally.

A CNN works by sliding (small) filter across the image, outputting the convolution (inner product) of the filter and pixels covered.



What is a CNN?

CNNs are a type of neural network, which works well with spatially dependent data (typically images). CNNs use parameter/weight sharing.

Multi-layered images (e.g. RGB or ATLAS calorimeter) are handled naturally.

A CNN works by sliding (small) filter across the image, outputting the convolution (inner product) of the filter and pixels covered.













