Applied ML Past Experiences with ML











"Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!"

On experience

Sentence:

"Experience is simply the name we give our mistakes",

[Oscar Wilde]

Lemma: *"I didn't fail. It was a learning experience",*

[Anonymous]

First encounters

Not having any experience with ML, I did a lot of mistakes:

- No description of architecture!
- No HP optimisation.
- No check of data-MC correspondence.

Worst of all, I had not thought of any way to cross check and calibrate the output.

But... simply throwing myself at it was a great experience to build on.



Higgs Search/Discovery

Motivation

Problem:

Given a number of clean ZZ events, determine if they are Higgs or SM diboson events!

Possible solution:

Since Higgses are produced quite differently then SM diboson ZZ, their angular distributions differ!



Variables available/used:

- Higgs rapidity
- Angle Z to Higgs in Higgs CM
- Angle lep- to Z in Z CM
- Angle lep- to Z* in Z* CM
- Fraction of mZ+mZ* to mHiggs



Generator level comparison



After fiducial requirements



Combining variables

Using the 5 variables (i.e. including rapidity) in a BDT (100 trees, 4 nodes):



BDTG response

Combined angular variable



Combined angular variable



TMVA BDT Response

Combined angular variable



Check for overtraining

Using 9 variables in a BDT (200 trees, 4 nodes) and checking for overtraining:

TMVA overtraining check for classifier: BDTG







Angular BDT score

197



Angular BDT score

19

Lessons Learned:

• Separation changed dramatically, when fiducial GeV) cuts were included. round • Very hard to include ML output in fits. • It is complicated to calculate systematics on ML output - one needs a plan (we didn't have one). It was nice to see, that there was no correlation between ML output and H candidate mass. But the results build confidence in our results in the Higgs to ZZ* group, and it subsequently became the path forward. -2 **ATLAS - Work in progress** -3 120 110 130 140 150 160

Angular BDT score

Higgs candidate mass (GeV)

Housing Prices

Individuel estimates

Shapley-values also gives the possibility to see the reason behind **individuel estimates**. Below is an example, illustrating this point.



Above is shown which factors that influences the final estimate of the sales price (and how much). The estimate is the sum of the contributions (here 6.86 MKr.).

This is a fantastic tool to get insight into the ML workings!!!

Word ranking



Result of including text

Natural Language Processing

Term Frequency - Inverse Document Frequency: TF-IDF

Natural weighting of words

CountVectorizer, TfidfVectorizer

Assign a weight to each word, according to its frequency of use. weight_IDF = log(N_{all} / N_{appearances})

MAD(XGB, numerics only) = 0.165

MAD(XGB, text only, BOW) = 0.254

MAD(XGB, combined) = 0.147

(Numerics: GeoPostNr, BeregnetAreal, ByggeAAr, EjendomsVaerdi0, Afstand_Kyst)

Result of including text

Lessons Learned:

- The ML part of the project was fun and BDTs worked really well.
 - Term F• Including text was (at the time) harder, but we
had a way to cross check, if it worked.-IDF
 - We were not at all prepared for the reluctance to use this in the real world.

"Big ships turn only slowly!"

each word, quency of use. N_{all} / N_{appearances})

MAD(XGB, numerics only) = 0.165

MAD(XGB, text only, BOW) = 0.254

MAD(XGB, combined) = 0.147

Coui

(Numerics: GeoPostNr, BeregnetAreal, ByggeAAr, EjendomsVaerdi0, Afstand_Kyst)

Electron Identification

Input Feature Ranking

Here is an example from particle physics. The blue variables were "known", but with SHAP we discovered three new quite good variables in data.



Input Feature Ranking

We could of course just add all variables, but want to stay simple, and training the models, we see that the three extra variables gives most of gain.



Input Feature Ranking



Electron Regression



ML at Work:



Electron Energy regression with CNN

Malte Algren^{*}, Aske Rosted, and **Troels C. Petersen** Niels Bohr Institute, Copenhagen (*now at Univ. of Geneva)

Outline

Outline of talk:

- Motivation
- Context
- Training a CNN for energy reconstruction:
 - The data
 - The selections
 - The input variables
 - The network architecture
 - Feature wIse Linear Modulation (FiLM)
- Results in MC
- Results in data (v1)
- Training in data and "simultaneous training"
 - Results in data (v2)
- Outlook



Motivation

Points of motivation:

- Improve $H \rightarrow ZZ^*$ and $H \rightarrow \gamma\gamma$ analyses
- Optimise searches for:
 - HH $\rightarrow \gamma \gamma bb$
 - $H \rightarrow Z\gamma$
 - $H \rightarrow \gamma^* \gamma$
- Improve resilience to pile-up
- Improve $Z \rightarrow$ ee reconstruction
- Utilise excellent data for testing:
 - CNN and GNN models
 - data+MC simultaneous training
 - $e+\gamma$ simultaneous training
- Improve non-Higgs searches

Goals of lecture:

- Give example of regression with CNN.
- Illustrate concept of attention and FiLM technique.
- Illustrate "target mismatch" and combined training.



Motivation

Points of motivation: (You don't have to care - just know the list is long!)

- Improve $H \rightarrow ZZ^*$ and $H \rightarrow \gamma\gamma$ analyses
- Optimise searches for:
 - HH $\rightarrow \gamma \gamma bb$
 - $-H \rightarrow Z\gamma$
- $H \rightarrow \gamma^* \gamma$
- Improve resilience to pile-up
 Improve Z → ee reconstruction
- Utilise excellent data for testing:
 - CNN and GNN models
 - data+MC simultaneous training
 - $e+\gamma$ simultaneous training
- Improve non-Higgs searches

Goals of lecture:

- Give example of regression with CNN.
- Illustrate concept of attention and FiLM technique.
- Illustrate "target mismatch" and combined training.







The scalars can be seen in table on the right.

The variables are both scalar and cell based.

Туре	Name	Description				
	Eacc	Energy deposit in layer 1-3 of ECAL.				
	η_{index}	η cell index of cluster of layer 2.				
	f0 _{cluster}	Ratio of energy between layer 0 and E_{acc} in $ \eta < 1.8$ (end of layer 0).				
Energy	R12	Ratio of energy between layer 1 and 2 in the ECAL.				
	p_t^{track}	p_T estimated from tracking for the particle (on e).				
	E_{TG3}	Ratio between the energy in the crack scintillators and E_{acc} within 1.4 < $ \eta $ < 1.6.				
	$E_{tile-gap}$	Sum of the energy deposited in the tile-gap.				
	η	Pseudorapidity of the particle.				
	$\Delta \phi_2^{rescaled}$	Difference between ϕ , as extrapolated by tracking, use for ECAL momentum estimation and ϕ of the ECAL cluster.				
	$\eta_{ m ModCalo}$	Relative η position w.r.t. the cell edge of layer 2 in the ECAL*.				
Geometric	$\Delta \eta_2$	Difference between η , as extrapolated by tracking, use for ECAL momentum estimation and η of the ECAL cluster (only <i>e</i>).				
	poscs ₂	Relative position of η within cell in layer 2 in ECAL. $2(\eta_{cluster} - \eta_{maxEcell})/0.025 - 1$, $\eta_{cluster}$ is η of the barycenter of the cluster and $\eta_{maxEcell}$ is η of the most energetic cell of the cluster.				
	$\Delta \phi_{TH3}$	Relative position in ϕ in a cell. mod $(2\pi + \phi, \pi/32) - \pi/32$.				
	$\langle \mu \rangle$	Average proton-proton interaction per bunch crossing.				
Misc.	n _{tracks}	# of tracks assigned (only e).				
	$n_{vertexReco}$	Number of reconstructed vertices.				

The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information: • Energy (primary variable)

• Time of cell energy



The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information:

- Energy (primary variable)
- Time of cell energy

In order to have the **same resolution** in each layer, we **upsample** the layers to the lowest common resolution (work by Lucas Erhke).





The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information:

- Energy (primary variable)
- Time of cell energy

Finally, we consider the (up to) 10 nearest tracks in a "TrackNet" input:

Туре	Name	Description					
Energy	p _{t,track} /q _{track}	Transverse momentum of track divided by its charge q					
	d_0/σ_{d0}	<i>d</i> 0 is the signed transverse distance between the point of closest approach and the <i>z</i> -axis where σ_{d0} is its uncertainty					
Geometric	ΔR	$\Delta R = \sqrt{(\phi_0 - \phi)^2 + (\eta_0 - \eta)^2}$					
	vertex _{track}	Reconstructed vertex of the track					
	<i>z</i> ₀	Longitudinal distance between the point of closest approach and the <i>z</i> -axis.					
	η_{track}	Reconstructed $ \eta $ of tracks.					
	ϕ_{track}	Reconstructed ϕ of tracks.					
	n_{pixel}	Number of hits in the pixel detector					
Misc.	n_{SCT}	Number of hits in the SCT					
	n _{TRT}	Number of hits in the TRT					



The network architecture

There are many ways to combine the input variables, and we have considered the following architectures, where the dashed lines are the considerations.



First, let us consider each part...

Feature wIse Linear Modulation



The network architecture

Testing all the different combinations yields the optimal architecture.

We evaluate the performance in the same way as previously done, namely the effective InterQuantile Range (eIQR) of the Relative Error (RE).

$eIQR = \frac{P_{75}(RE) - P_{25}}{1.349}(RE)$, $RE = \frac{E_{calib}}{E_{calib}}$,			-	Hyperparameter	Parameter
				Ltru	ıth		Units	(128, 64, 32, 16)
							Normalization	Batch
F							Kernel size & filters	5
		reIÇ	2R75	reIQR95			Connected to	[Top]
3	Basic	-(0.121	-0.025	E.S.		ScalarNet	
1	FiLM: scalar FiLM: scalar - top: scalar		0.229	0.257 0.252	10		Units	(256)
			0.220				Normalization	Batch
	Fil M: scalar - top: scalar track	0.222		0.251	2	2	Connected to	[FiLM]
ſ	FILM. Scalar - top. Scalar track		5.223	0.251	Dect A.		FiLM gen.	
	FiLM: scalar - top: track		0.226	0.264	Dest	A	Contrecture	(512, 1024)
	FiLM: scalar track	(0.228	0.265	3	5.00	Normalization	Batch
	FiLM: scalar track - top: scalar track FiLM: track - top: scalar		0 .2 10	0.262	2		CNNnet	
25			0.042	-0.067	3		Down-sampling	MaxPool
	Fil M: track top: track	0.		0.007	87		Globalpooling	MaxPool
	top: scalar top: scalar track		5.140	0.149			Number of blocks	3
			0.154	4 -0.131 3 0.233			Depth of blocks	4
			0.213				Тор	
	top: track	(0.136	0.164			Units	(512, 512, 1)
					1		Output activation	ReLU



The results in 2D - MC

The E_T distribution for truth (x-axis) and reconstruction (y-axis) can be compared for the current ATLAS and the DeepCalo algorithms.

As the figure shows, both algorithms do well, and improve with energy.

As the statistics is largest around 40 GeV, this is where the comparison is most detailed, and here DeepCalo visibly has a significantly reduced lower edge. Thus, the DeepCalo more rarely undershoots the energy.



The results in 1D - MC

Integrating the previous plot into 1D considering the RE distribution, we see a general sharpening. The improvement in relative eIQR (reIQR) is about 22%.



Naively, we would of course love to see a similar number in data!

Result in Zee - MC

On the Zee peak, we evaluate the improvement by fitting with a BW⊗CB fit, considering the CB width (sigmaCB) as the performance parameter. We get:



Result in Zee - MC

On the Zee peak, we evaluate the improvement by fitting with a BW⊗CB fit, considering the CB width (sigmaCB) as the performance parameter. We get:



Results on Zee - data (v1)

The result we get is a much more modest improvement:

$$\langle 1 - rac{\sigma^{DeepCalo}_{CB}}{\sigma^{ATLAS}_{CB}}
angle = 1 - rac{2.058 \pm 0.010}{2.271 \pm 0.019} = 9.4 \pm 0.9\%.$$

Though perhaps a little disappointing, this is not surprising, as we can not expect the MC to mimic data perfectly in the very large space considered. Also, models trained on Zee do not generalise well to all energies (EG, 6.8%).



45





Training in data

Using Zee events with invariant masses 86-97 GeV, one can get "approximate labels" in data, by assuming the true Z mass: $M^2 = 2p_{T,1}p_{T,2}(\cosh(\eta_1 - \eta_2) - \cos(\phi_1 - \phi_2)), p_T$

Using such labels, we train in data and get...

$$M^{2} = 2p_{T,1}p_{T,2}(\cosh(\eta_{1} - \eta_{2}) - \cos(\phi_{1} - \phi_{2})), \quad p_{T} = E_{T} \updownarrow$$
$$E_{label,data} = \frac{M^{2}}{2E_{T,2}(\cosh(\eta_{1} - \eta_{2}) - \cos(\phi_{1} - \phi_{2}))'}$$
with $E_{T,2} = E$ calib^(BDT) and $M^{2} = 91.19^{2}$





Training in data

Using Zee events with invariant masses 86-97 GeV, one can get "approximate labels" in data, by assuming the true Z mass: $M^2 = 2p_{T,1}p_{T,2}(\cosh(\eta_1 - \eta_2) - \cos(\phi_1 - \phi_2)), p_1$

Using such labels, we train in data and get...







Training in data and MC

Once we have labels in data, there is nothing keeping us from combining the loss functions of MC and data (they even have the same form), and thus training **simultaneously** in data and MC:

$$\mathcal{L}(y, \hat{y}) = \mathcal{L}(y_{(\text{Zee, MC})}, \hat{y}_{(\text{Zee, MC})}) + \mathcal{L}(y_{(\text{Zee, Data})}, \hat{y}_{(\text{Zee, Data})})$$

This allows the model to both use the "strength" of MC, but also learn the differences between MC and real data.

Doing this and trying out the result in MC first yields:

$$\langle reIQR_{75}^{DeepCalo}
angle = 22.1 \pm 0.3\%$$

OK, so at least it doesn't ruin the model for MC. Now let us try data...



Result in data (v2)

The result in data is rather encouraging, and **greater than the sum of the improvements** from training separately in MC (9.4%) and data (5.9%).



Outlook

While this is still "only" an improvement in the electron energy regression, and only for lower energies (Zee range), the simultaneous training allows for extending the energy range, by including the Electron Gun MC.

Furthermore, this training might be extended to include photons, as these behave much the same as electrons, and suffer the same sources of uncertainties and smearing.

For improving the H $\rightarrow \gamma \gamma$ resolution, one might use the following loss function and related training samples:

$$\begin{aligned} \mathcal{L}(y, \hat{y}) &= \mathcal{L}(y_{(\text{Zee, MC})}, \hat{y}_{(\text{Zee, MC})}) + \mathcal{L}(y_{(\text{Zee, Data})}, \hat{y}_{(\text{Zee, Data})}) + \\ \mathcal{L}(y_{(Z\mu\mu\gamma, \text{MC})}, \hat{y}_{(Z\mu\mu\gamma, \text{MC})}) + \mathcal{L}(y_{(Z\mu\mu\gamma, \text{Data})}, \hat{y}_{(Z\mu\mu\gamma, \text{Data})}) + \\ \mathcal{L}(y_{(H\gamma\gamma, \text{MC})}, \hat{y}_{(H\gamma\gamma, \text{MC})}) \end{aligned}$$

Meanwhile, we are trying to write this up somehow (but Malte is now a Ph.D. in Geneva).

Outlook

While this is st only for lower extending the

Furthermore, t behave much t and smearing. Lessons Learned:

- Remember to think about publishing. Even what may seem "a fun little example" at the time, may turn out to inspire a new line of thinking.
- Remember to think about the longevity of any approach. In this case, the storage of cell information was discontinued shortly after!

gression, and ws for

as these funcertainties

For improving the resolution, one might use the following loss function and related training samples:

 $\begin{aligned} \mathcal{L}(y, \hat{y}) &= \mathcal{L}(y_{(\text{Zee, MC})}, \hat{y}_{(\text{Zee, MC})}) + \mathcal{L}(y_{(\text{Zee, Data})}, \hat{y}_{(\text{Zee, Data})}) + \\ \mathcal{L}(y_{(\text{Z}\mu\mu\gamma, \text{MC})}, \hat{y}_{(\text{Z}\mu\mu\gamma, \text{MC})}) + \mathcal{L}(y_{(\text{Z}\mu\mu\gamma, \text{Data})}, \hat{y}_{(\text{Z}\mu\mu\gamma, \text{Data})}) + \\ \mathcal{L}(y_{(H\gamma\gamma, \text{MC})}, \hat{y}_{(H\gamma\gamma, \text{MC})}) \end{aligned}$

Meanwhile, we are trying to write this up somehow (but Malte is now a Ph.D. in Geneva).

DeepFRET

FRET is a technique used to study and dynamics of biomolecules. The data is a "trace", which is a time series with possible phase transitions.

The group would go through 10000 traces and select about 250 of these... **by hand**!!! This took a few people about a week, and was neither reproducible nor optimal.

So we made DeepFRET.



DeepFRET



So we made DeepFRET.

Knee- & Hip surgery

ROC curve

The previous figure is summarised in this plot, where one can see the false positive rate (x-axis) vs. the true positive rate (y-axis).



The red dot corresponds to the cut before (> 0.14), and yields the values shown on the right.

From a medical point of view, one can then choose an operational point on the blue curve (the dashed line being a random choice).

Ranking of features



mean([SHAP value]) (average impact on model output magnitude)

Further improvements

We don't know which is "Hospital=9", but we don't want to send Mathias there!



Further improvements

We don't know which is "Hospital=9", but we don't want to send Mathias there!

			2.00
		Lessons Learned:	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	0.4	• The enquiry about the data was fitting, and in	
	0.4 -	this case the data was really nice.	
		• BDTs were the obvious way to go, given all	1.25
L	0.2 -	sorts of NaNs, categories, and binary input.	
e fo		• The speed with which we could make models	SUS
aluo oita		impressed our collaborators - twice.	tat
sv c	• Asking for outline data is useful.	0.50 m	
HAI L	• The use of SHAP values was extremely useful.	<u><</u> .	
S	1. m	and also convinced our colleagues.	
	the second		-0.25
6	-0.2 -		
	Same		
0	-0.4 -	1 2 3 4 5 6 7 8 9	- 1.00
		hospital	

IceBoost

Estimating the volume of glaciers is "hard" given the lacking 3D view. But it can be done using satellite images, climate, and physics (mass balance).



In order to estimate what the ground underneath looks like, we tried using inpainting. It worked reasonably well, but never beat the BDT approach.

IceBoost

Estimating the volume of glaciers is "hard" given the lacking 3D view. But it can be done using satellite images, climate, and physics (mass balance).



In order to estimate what the ground underneath looks like, we tried using inpainting. It worked reasonably well, but never beat the BDT approach.