### Higgs Machine Learning Challenge what now?



David Rousseau LAL-Orsay

rousseau@lal.in2p3.fr

17<sup>th</sup> Nov 2015, DESY

### ... in a nutshell

Why not put some ATLAS simulated data on the web and ask data scientists to find the best machine learning algorithm (=MVA) to find the Higgs ?

- Instead of HEP people browsing machine learning papers, coding or downloading possibly interesting algorithm, trying and seeing whether it can work for our problems
- Challenge for us : make a full ATLAS Higgs analysis simple for non physicists, but not too simple so that it remains useful
- Also try to foster long term collaborations between HEP and ML

David Rousseau, HiggsML what now, 16th November 2015

### Dataset

Permanently available and usable by anyone (also Primitive 3-vectors allowing to compute the conf non ATLAS) on CERN Open Data: note variables (mass neglected), http://opendata.cern.ch/collection/ATLAS-Higgs-Challenge-2014 16 independent variables: ASCII csv file, with mixture of Higgs to tautau PRI tau pt (lephad) signal and corresponding backgrounds, PRI tau eta from official GFANT4 ATLAS simulation PRI tau phi Weight and signal/background label (for training dataset only) PRI lep pt weight (fully normalised) PRI lep eta label : « s » or « b » PRI lep phi Conf note variables used for categorization or BDT: PRI met DER mass MMC PRI met phi DER mass transverse met lep PRI met sumet DER\_mass\_vis PRI\_jet\_num (0,1,2,3, capped at 3) DER pt h PRI jet leading pt DER deltaeta jet jet PRI jet leading eta VBF VBF DER\_mass\_jet\_jet PRI jet leading phi signature signature DER prodeta jet jet PRI\_jet\_subleading\_pt DER deltar tau lep PRI jet subleading eta DER pt tot PRI jet subleading phi DER sum pt PRI jet all pt DER\_pt\_ratio\_lep\_tau DER met phi centrality HiggsML what now, 16th November 2015 David Rousseau. 3 DER lep eta centrality

# How did it work ?

- □ First idea in Sep 2012
- Challenge ran from May to September 2014
- People register to Kaggle web site hosted <u>https://www.kaggle.com/c/higgs-boson</u>. (additional info on <u>https://higgsml.lal.in2p3.fr</u>)
- Open to almost any one
  - o Data scientist
  - HEP physicists
  - o Students, geeks,
  - Except LAL-Orsay employees (for legal reasons)
- ...download training dataset (with label) with 250k events
- ...train their own algorithm to optimise the significance (à la s/sqrt(b))
- ...download test dataset (without labels) with 550k events
- ...upload their own classification
- The site automatically calculates significance. Public (100k events) and private (450k events) leader boards update instantly. (Only the public is visible)
- Competition closed mid september 2014. Private leaderboard is disclosed. People are asked to provide their code and methods. Best 1 2 3 win 7k\$ 4k\$ 2k\$
- □ In addition, the potentially most interesting one gets the "HEP meets ML award"

#### Funded by: Paris Saclay Center for Data Science, Google, INRIA



# Real analysis vs challenge

- 1. Systematics (and data vs MC)
- 2. 2 categories x n BDT score bins
- 3. Background estimated from data (embedded, anti tau, control region) and some MC
- Weights include all corrections.
  Some negative weights (tt)
- 5. Potentially use any information from all 2012 data and MC events
- 6. Few variables fed in two BDT
- 7. Significance from complete fit with NP etc...
- 8. MVA with TMVA BDT

- 1. No systematics
- 2. No categories, one signal region
- 3. Straight use of ATLAS G4 MC
- Weights only include normalisation and pythia weight. Neg. weight events rejected.
- 5. Only use variables and events preselected by the real analysis
- All BDT variables + categorisation variables + primitives 3-vector
- 7. Significance from "regularised Asimov"
- 8. MVA "no-limit"

#### Simpler, but not too simple!

David Rousseau, HiggsML what now, 16th November 2015

# Significance

- Need to have one robust estimator of the quality of the classification algorithm
- Decided to use the well known (in HEP) "Asimov" formula (G. Cowan, K. Cranmer, E. Gross, and O. Vitells, "Asymptotic formulae for likelihood-based tests of new physics", *EPJCC*, vol. 71, pp. 1–19, 2011.) with regularization on top
  - o  $\sqrt{(2^*((s+b')^*\log(1+s/b')-s))}$  ~s/√b'
  - with s and b'=b+10 normalised to 2012 data taking luminosity:
  - $s=\Sigma$ (selected signal) weights\_i
  - $b=\Sigma$ (selected background) weights\_i
- Why b'=b+10 ("regularisation") : practical way to avoid large significance fluctuation when small phase space region with very few background events is chosen. Do not want to pick winners on their luck.
- Note that normalisation already included in the weights : no need to explain integrated luminosity and cross-section

7

X

(C)

## What data did we release ?

□ From ATLAS full sim Geant4 MC12 production

- 30 variables
- □ Signal is H→tautau, Background a mixture of : Z, top, W
- Based on November 2013 ATLAS Htautau conf note ATLAS-CONF-2013-108
- Preselection for lep-had topology : single lepton trigger, one lepton identified, one hadronic tau identified
- □ →800.000 events:
  - o 250.000 training data set
  - 550.000 test data set without label and weight
- Reproduces reasonably well (~20%) content of 3 highest sensitivity bins (x 2 categories) in conf note
- (some background and many correction factors deliberately omitted so that the sample cannot be used for physics, only for machine learning studies)

# **Participation**

- Big success !
- 1785 teams (1942 people) have participated (participation=submission of at least one solution)
  - o (6517 people have downloaded the data)
  - →most popular challenge on the Kaggle platform (until spring 2015)
  - o 35772 solutions uploaded
- 136 forum topics with 1100 posts
- Many participants have worked very hard

# **Final leaderboard**

#	∆rank	k Team Name ‡ model uploaded * in the money		Score 😨	Entries	Last Submission UTC (Best – Last Submission)
1	<b>↑1</b>	Gábor Melis ‡ *	7000\$	3.80581	110	Sun, 14 Sep 2014 09:10:04 (-0h)
2	<b>↑1</b>	Tim Salimans ‡ *	4000\$	3.78913	57	Mon, 15 Sep 2014 23:49:02 (-40.6d)
3	<b>↑1</b>	nhlx5haze ‡ *	2000\$	3.78682	254	Mon, 15 Sep 2014 16:50:01 (-76.3d)
4	<b>↑38</b>	ChoKo Team 🎤		3.77526	216	Mon, 15 Sep 2014 15:21:36 (-42.1h)
5	<b>↑35</b>	cheng chen		3.77384	21	Mon, 15 Sep 2014 23:29:29 (-0h)
6	<b>↑16</b>	quantify		3.77086	8	Mon, 15 Sep 2014 16:12:48 (-7.3h)
7	<b>↑1</b>	Stanislav Semeno	ov & Co (HSE Yandex)	3.76211	68	Mon, 15 Sep 2014 20:19:03
8	↓ <b>7</b>	Luboš Motl's tear	n #	3.76050	589	Mon, 15 Sep 2014 08:38:49 (-1.6h)
9	<b>↑8</b>	Roberto-UCIIIM		3.75864	292	Mon, 15 Sep 2014 23:44:42 (-44d)
10	<b>↑2</b>	Davut & Josef 🎤		3.75838	161	Mon, 15 Sep 2014 23:24:32 (-4.5d)
45	<b>↑5</b>	crowwork 🌶 ‡	HEP meets ML award XGBoost authors Free trip to CERN	3.71885	94	Mon, 15 Sep 2014 23:45:00 (-5.1d)
782	2 ↓149	Eckhard	Tuned TMVA	3.4994	5 29	Mon, 15 Sep 2014 07:26:13 (-46.1h)
99	1 <b>↑4</b>	Rem.		3.20423	2	Mon, 16 Jun 2014 21:53:43 (-30.4h)
8	4	simple TMVA b	oosted trees	3.19956		

### **Best private scores**



David Rousseau, HiggsML what now, 16th November 2015



David Rousseau, HiggsML what now, 16th November 2015

### **TMVA vs Gabor**





- vbf, boosted categories as is ATLAS note (no ATLAS insider information)
- tmva, gabor are trained without categories, on full 30 variables (not directly comparable to ATLAS analysis)
- (also significance is simple asimov, no bin, no systematics (and fake tau missing))
- Gabor improves more significantly in VBF categories (2 jets →events more complex)

David Rousseau HiggsML visits CERN, 19th May 2015

### Lessons



### What did we learn

Very successful full day satellite workshop at NIPS (one of the two major Machine Learning conferences) in Dec 2014 @ Montreal:

https://indico.lal.in2p3.fr/event/2632/

- Proceedings just published (August 2015 : JMLR Workshop and Proceedings Vol 42 <u>http://jmlr.org/proceedings/papers/v42/</u>) Contributions from some of the top players, plus summary from organisers
- Many additional piece of information in kaggle forum or random blog or github repository
- □ Each participant have used a range of ideas selected by trial and error → difficult to decipher what really worked best at the end

# Imputation

In ML jargon, this is the handling of missing variables, a very hot topic

In HiggsML, we provided leading and sub-leading jet 4momenta, plus variable based on these (e.g. di-jet mass), but many events with just one or zero jet

• In addition MMC would fail in a few percent of the cases

- □ No clear winning strategy among:
  - o not doing anything special
  - Replace missing variables by average on other events
  - Separate training samples according to available variables

# **Algorithms**

- "deep" Neural Nets win (Gabor Melis).
  - Gabor's words : "deep" in 2014 because using 3 hidden layers, would not qualified as deep nowadays
- BDT marginally behind (number 2 was 0.02 behind in significance)
  - Gabor's words : NN not worth it, too much work/too many possibilities to tune the training
  - (Gabor has just been hired by DeepMind)
- Meta-ensemble (combining BDT or NN with different hyper parameters) marginally better, but much more complex, not worth it
- Conclusion : for a typical HEP problem, BDT should be the default choice (OK we sort of knew about it)

### Software

- Lots of development of Machine Learning Open Source software outside HEP
- □ In particular:
  - XGBoost (eXtreme Gradient Boosting): released for the HiggsML challenge, used by many participants. Now used in other challenges as well. One of the best software on the market for BDT/BRT. Good performance out of the box. Also fast (multithreaded)
  - SciKit-learn : large developer/user base. Toolbox like TMVA. Used already a bit in ATLAS (e.g. Htautau hadhad channel at least)
  - Note that both software were improved thanks to the challenge, in particular to handle event weights
  - Note that quite often these Open Source software are routinely multithreaded, and sometimes even run on GPU (NN, not BDT)

#### TMVA:

- New effort to rejuvenate TMVA on three fronts (see Saas-Fe Root user workshop Sep 2015)
- Improve TMVA algorithms
- Improve TMVA structure and user interface (i.e. CV, see later)
- Interface to the outside world (e.g. R interface functional now)

# **Feature Engineering**

- In ML jargon, this is the building of new variables from the original ones
- We (HEP) have been doing this since the beginning of times
- Given enough training data, ML techniques could "discover" these features (e.g. invent the concept of transverse mass)
- □ It did not work at all for HiggsML (significance less than 3 (wrt 3.8) if removing the high level variables provided
- There are techniques to automatically generate new features
- Not clear they would beat HEP expertise

## **Cross Validation**

- Cross Validation (CV) are techniques to measure MVA performance independently of the training
- Goal is to build an optimisation curve (e.g. significance, ROC,..) with the smallest variance (despite lack of data), for a better optimisation of hyper parameters or choice of techniques
- Default TMVA CV (one fold CV):
  - split sample in two halves A and B.
  - o train on A, test on B
- Two-fold CV (e.g. ATLAS Htautau analysis)
  - o Split sample in two halves A and B
  - Train on A, test on B; train on B test A
  - →test statistics = total statistics→double test statistics wrt one fold CV (double training time of course)
- □ Five-fold CV (e.g. Gabor)
  - Split sample in 5 equal pieces A,B,C,D and E
  - Train on ABCD, test on E;train on ABCE, test on D; etc...
  - → same test statistics wrt Two-fold CV, but larger training statistics 4/5 over ½ (larger training time as well)
- CV à la Gabor (he did not invent it but no better name)
  - o Redo Five-fold CV e.g. 4 times with a different random splitting
  - For each event, one has now 4 different (but similar) scores. Then
    - Average these scores (with whatever definition of "average")
    - Or build directly the optimisation curve from the 4\*N scores, with additional weight ¼→smoother curve



# **Focussed learning**

- By default, classification algorithm will optimise the overall ROC curve (typically the Area Under roc Curve, AUC)
- However we often have more specific figure of merit, like the significance à la s/sqrt(b) (which depend of the size of the subsample) (not to mention full blown RooFit !)



AN

- ❑ We want to focus on a specific region of the ROC curve (not AUC), with background rejection >95% and signal efficiency ~10-20%
- Different techniques have been used by participants to handle this:
  - Hyper parameter optimisation maximising significance (but quickly overtraining)
  - Chose internal training parameter as a function of the optimisation functional
  - Prescription to modify the event weights and iterate learning (Weighted Classification Cascade, see arxiv 1409.2655)
  - $\rightarrow$  need more study to understand what works best
- Note : systematics were deliberately ignored for the challenge, but these methods could be used with a significance including systematics



# Handling systematics

- □ Typically we would (i) train our BDT (ii) compute the systematics
- ➡ how to tell the BDT to avoid poorly controlled variable or background?
- Open problem. Can be tackled e.g. with focussed training.
- New topic for Machine Learning. They are quite excited about it. E.g. :



David Rousseau, HiggsML what now, 16th November 2015

### **Random additional ideas**

The following were also mentioned in brainstorming at the NIPS workshop or in other discussions triggered by the challenge (not complete!)

- Use of deep learning (see papers by Baldi, Sadowski, Whiteson): able to guess high level feature on some problem but not others ?
- Handle the lack of training statistics (in particular for deep learning) by combining in a clever way full/fast sim (almost accurate but expensive) and super fast sim (~Delphes) (inaccurate but cheap). ML jargon "Transfer Learning"
- Transform a classification problem in a regression problem, easier to train, with a surrogate function. E.g. build a very sophisticated deep NN, train it, then emulate it with a simple BDT
  - The simple BDT does not work better than the very sophisticated deep NN, but it is faster and better than training directly the simple BDT

### **Other challenges**



# **LHCb : flavour of physics**

- □ LHCb organised in summer 2015 another challenge "flavour of physics": search for LFV decay  $\tau \rightarrow \mu \mu \mu$
- similar to HiggsML, with a big novelty:

- o some variables known to be poorly described by MC
- algorithm had to behave similarly on data and MC in a control region  $D0 \rightarrow K\pi\pi$
- ➡Nice idea, however, never underestimates the machine learners: They devised an algorithm which
  - was able to distinguish control region from signal region
  - was behaving well (data=MC) in the control region
  - but was recklessly abusing the data/MC difference in the signal region
- □ → rules had to be changed in the middle of the challenge to disallow this

#### Anyway, this does show that systematics is tricky to handle



# **Tracking with pileup**

#### Graeme Stewart ECFA HL-LHC workshop 2014

- Tracking 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)









# **Pattern recognition**

- Pattern recognition, connecting the dots, is a very old, very hot topic in Artificial Intelligence
- □ Just one example among many from NIPS 2014 : <u>http://papers.nips.cc/paper/5572-a-complete-variational-tracker.pdf</u>



- Note that these are realtime applications, with CPU constraints
- Worry about efficiency, "track swap" avid Rousseau HiggsML and tracking challenges CTD 2015 Berkeley



# **Tracking challenge ?**

- Trickier to organise than HiggsML or the like:
  - less "on-the-shelf" algorithms than for classification
  - Figure of merit combination of efficiency/fake rate/CPU time
  - CPU time to be measured in a well defined way
- Goal is to go online in summer 2016



### Summary

- Wealth of disorganised input from the challenge participants. What we could decipher:
  - BDT still the algorithm of choice
  - Better software out there (XGBoost, SciKitLearn) than *current* TMVA (but TMVA re-boost effort started)
  - Many techniques beyond just BDT training (Cross Validation, focussed training etc...)
  - Lots of expertise in ML community we should tap into

- Pointer collection:
  - o <u>https://www.kaggle.com/c/higgs-boson</u>
  - o <u>https://higgsml.lal.in2p3.fr</u>
  - <u>http://opendata.cern.ch/collection/ATLAS-Higgs-Challenge-2014</u>: permanent home of the challenge dataset
  - <u>https://indico.lal.in2p3.fr/event/2632/</u> NIPS 2014 workshop agenda and NEW proceedings <u>http://jmlr.org/proceedings/papers/v42/</u>
  - o http://cern.ch/higgsml-visit mini workshop at CERN
- Mailing list just opened to any one with an interest in both Data Science and High Energy Physics (mainly for announcements) : <u>HEP-data-science@googlegroups.com</u>

# Outlook

- Initiatives are now sprouting:
- Challenges as mentioned
- Inter-Experimental LHC Machine Learning Working Group : <u>http://iml.cern.ch</u> and mailing list <u>lhc-machinelearning-wg@cern.ch</u> (initially for TMVA reboost, but larger scope)
- Workshop "Data Science @ LHC" at CERN, 9<sup>th</sup>-13<sup>th</sup> Nov 2015 <u>http://cern.ch/DataScienceLHC2015</u> (videos)

 Opportunities beyond BDT and Deep NN (e.g. Approximate Bayesian Processes, Gaussian Processes, etc...)

Last word : if you in HEP want to embark in this, there is probably a friendly Machine Learner next door at your home institute. Teaming with such people is awesome!

David Rousseau, HiggsML what now, 16th November 2015

