

The Higgs Boson Machine Learning Challenge

ATLAS Collaboration



HEP

Machine Learning

Doing Big Data for 30 years

Using MVA for 20 years (neural net, boosted decision trees, and more)

Large players in industry, and academics

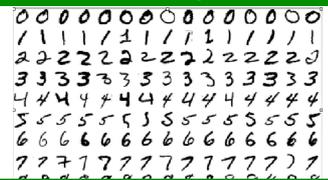
Large computing power

Complex problems

Exponential growth of the field. New algorithms and techniques

HiggsML challenge

Character recognition



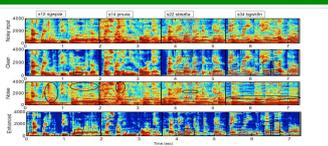
Real time face detection



Emotion recognition



Speech recognition



Let's put some ATLAS simulated data on the web and ask data scientists to develop the best machine learning algorithm to find the Higgs !

Jointly organised by ATLAS physicists and data scientists:

How does it work ?

- Participant register to kaggle
- Download training dataset (250 kEvt) with weight and sig or bkg label
- Train their best algorithm, adapt it to the specificities of the Challenge
- Download test dataset (550 kEvt) without weight and label
- Apply trained algorithm to test dataset, determine s/b classification

Public leaderboard from 100 kEvt

Rank	Participant	Score
1
2
3

Private leaderboard from 450 kEvt

Rank	Participant	Score
1
2
3

Best scores



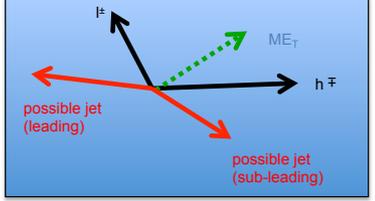
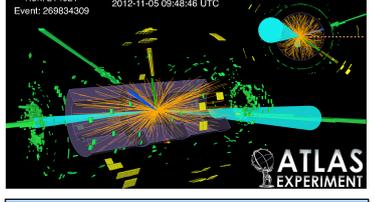
Best method



Choice of ATLAS H to tau tau analysis

- "Evidence for Higgs Boson Decays to tau tau Final State with the ATLAS detector" ATLAS CONF-2013-108
- 4.1 sigma observed (3.2 sigma expected)
- Direct evidence of Higgs coupling to leptons
- Complex analysis (e.g. signal H mass ~ one sigma from dominant Z to tau tau background, VBF signature...)

lepton-hadron topology



Choice of figure of merit

Need one and only one robust estimator of the quality of the classification algorithm : Approximate Medium Significance **AMS**

Given s and b expected number of signal and background normalised to 2012 luminosity:
 $s = \sum(\text{selected signal})$ weights,
 $b = \sum(\text{selected bkg})$ weights,

Decided to use the "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 » :
 $b = b + 10$ (avoid large fluctuations in small regions)

$$AMS = \sqrt{2} \sqrt{(s+b') \log(1+s/b') - s}$$

Real analysis vs Challenge

- Systematics
- 2 categories x n BDT score bins
- Background estimated from data (embedded, anti tau, control region) and some MC
- Weights include all corrections. Some negative weights.
- Potentially use any information from all 2012 data and MonteCarlo events
- Few variables fed in two BDT
- Significance from complete fit with Nuisance Parameters, Control Regions, etc...
- MVA with TMVA BDT

- No systematics
- No categories, one signal region
- Straight use of ATLAS Geant4 MonteCarlo : signal Higgs, backgrounds Z, W, top
- Weights include normalisation and generator weight. Negative weight events rejected.
- Only use variables and events preselected by the real analysis
- All BDT variables + categorisation variables + primitives 3-vector
- Significance from "regularised Asimov"
- MVA "no-limit"

Variables provided

Weight and signal/background label

weight
label

Conference note DER ived variables used for categorization or BDT (VBF and Boosted categories):

- | | |
|------------------------|-----------------------|
| DER_mass_MMC | DER_deltar_tau_lep |
| DER_mass_trans_met_lep | DER_pt_tot |
| DER_mass_vis | DER_sum_pt |
| DER_pt_h | DER_pt_ratio_lep_tau |
| DER_deltaeta_jet_jet | DER_met_phi_centralty |
| DER_mass_jet_jet | DER_lep_eta_centralty |
| DER_prodelta_jet_jet | |

PRImitive 3-vectors allowing to compute the DER variables (mass neglected)

- | | |
|---------------|------------------------|
| PRI_tau_pt | PRI_jet_num |
| PRI_tau_eta | PRI_jet_leading_pt |
| PRI_tau_phi | PRI_jet_leading_eta |
| PRI_lep_pt | PRI_jet_leading_phi |
| PRI_lep_eta | PRI_jet_subleading_pt |
| PRI_lep_phi | PRI_jet_subleading_eta |
| PRI_met | PRI_jet_subleading_phi |
| PRI_met_phi | PRI_jet_subleading_pt |
| PRI_met_sumet | PRI_jet_all_pt |

Simpler, but not simple!

<https://www.kaggle.com/c/higgs-boson>

Contact : higgsml@lal.in2p3.fr

