



### Machine learning methods to improve boosted Higgs boson tagging at ATLAS

Speaker: Vladlen G. Supervisor: Inês Ochoa

### A brief intro...

Why?

-To Study Beyond the Standard Model physics :

-There, heavy resonances may appear, which can decay into Higgs bosons with high pt (boosted).

-Reconstruction techniques for boosted higgs boson play an important role.







## Boosted Higgs and difficulties

b-jet





What are jets?

What's the difference between boosted and not boosted regimes?

Why is it so difficult to characterize them?



### b-tagging method What is b-tagging?

- b Quarks fragment into B-hadrons
- Displaced vertex (secondary vertex) from primary vertex due to its long life (~1.5ps)







## A deep NN and how to feed it

#### Architecture

Layers: **5 hidden** Nodes: **112, 96, 48, 24, 12, 6, 3** Optimizer: **Adam** Loss: **Categorical Crossentropy** Learning rate: **0.01** Batch size: **256** Activation function: **Relu** 

-Classes: Higgs, QCD (dijets), Top

Food: Data (treated)

With sauce: Features



ex: DL1rT\_pu\_1: probability distribution of the jet 1 being created from an up quark

**11 Features in Total!** 

ATLAS

7 FXPFRIMENT











Probability of three VR track-jets to be light, for large-R jets in Higgs,Top and QCD



8



### Results to reproduce...

https://cds.cern.ch/record/2724739/files/ATL-PHYS-PUB-2020-019.pdf







# We finally feed our NN!



### ROC curves



ROC with Y\_pred and Y\_test, for Classes: H,QCD,Top with 11 features



True Positive Rate vs False Positive Rate?

Ideal performance of the architecture?

Area is an indicator of the model's ability to distinguish between classes.

### **Confusion Matrices**





Allows the visualization of the performance of the algorithm, by comparing the jet class predicted by the NN and the real one.

### Discrimination variable



-Used to discriminate between classes -f\_{top}= 0.25 pHiggs, ptop,pmultijet predicted output (probability added to 1) from the NN, for the 3 classes.









### We do the same, but better...



#### 100 Epochs with 11 features:



## ROC (and roll) with some Confusion





-Top the most interchangeable -Most complex (it has W's) and it intersects the range of QCD and Higgs

### DXbb Distribution for 100 epoch





-The Intersection is visible here -We could also study the performance vs flavour composition of the jets in the simulation.

## Substructure Variables (adding 🖕





Adding to the information given from the b-tagging, we also could have additional knowledge of the structure of the jets via the calorimeter and the tracks.

This could be used to characterize better our system, among Higgs, QCD and Top



Adding a total of 14 additional variables for example "split12" that uses symmetry to distinguish asymmetric QCD splittings from heavy symmetric particle decays

## Now 100 epochs with 25 features



Loss / Mean Squared Error



### A new ROC and total CM







### Comparing Confusion matrices







### Dxbb distribution with 25 features









## ROC with signal eff. and background

With mass cuts in the range of [75, 145] GeV to define Higgs boson mass region



### Mass distribution, after Dxbb>1.8



With Dxbb>1.8, Mass distribution with 11 features

With Dxbb>1.8, Mass distribution with 25 features



## Something interesting appears...





With Dxbb>1.8, Mass distribution with 25 features (Normalized)



### **Conclusions:**



- Looking at the Confusion matrices, we identified that the **Top background was the most difficult to filter through** by being the most complex and by intersecting the ranges of Higgs and QCD.
- With the b-tagging variables alone, the **background does not peak in the Higgs mass region** after the cut in the DXbb distribution.
- The comparison between 11 features and 11+ substructure variables shows an increase in background rejection up to factor of 2.
- However with the substructure we can also distort the background jet mass distribution. As a next step this could also be corrected via a NN technique.

