# Jet flavor tagging with Deep Learning using Python

## 2nd Developers@CERN Forum

**Marie Lanfermann**, Tobias Golling

marie.lanfermann@unige.ch

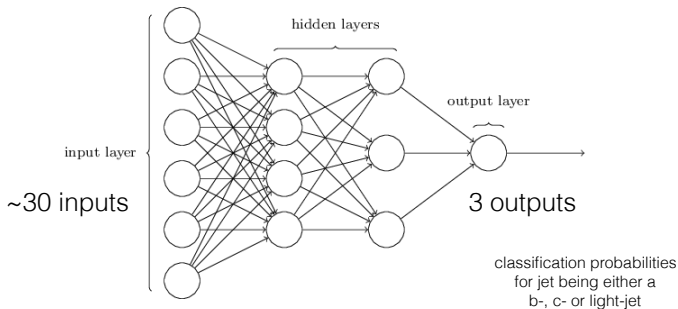University of Geneva

31st May 2016

## Outline

# Approach

Start: MC sample



End: NN configuration for classification problem

# Data preparation - part 1



ROOT Data Analysis Framework

flat ntuple

https://github.com/rootpy/root_numpy

**root_numpy**

NumPy

arrays

pandas

DataFrame

*PrepareSamples.py*

apply basic selection cuts

HDF

| step | resulting change |
|------|------------------|
| determine training and testing set separation (variable via argument) | + new pandas.DataFrame storing index lists |
| 2D (eta, pT) reweighting (distribution shape via argument) | + new pandas.Series to add to DataFrame |
| default variable resetting | + new pandas.Series check variables & pandas.DataFrame update predefined or calculated on valid values |

changes performed using functions defined in python module *btag_nn*

*btag_nn_inputs.py*

sample-specific definitions:
• binning
• input variable details (lists of strings) incl. calculation of new defaults

prepared pandas.DataFrames:

HDF

• preprocessed data
• index lists for training, testing
• default dictionary

python

# Data preparation - part 2

Normalize NN
input data:
calculate & apply
scale and offset

**HDF** & **JSON**
data-details

numpy arrays
for reloading in next run with identical
hyperparameter-independent basic setting

contains for each NN input variable:
• order
• default value
• input normalization (scale and offset)

```
>>> df.as_matrix()
array([[-0.13711547, -0.95178095,  2.44100016],
       [ 0.68082075,  0.56670109,  0.32898284],
       [ 0.66057529,  1.34594355,  0.35998082],
       [ 0.03709011, -1.14934791, -1.80300476],
       [ 0.60576849,  0.04855424,  0.60367536],
       [ 0.97191416, -0.43893764,  0.33687459],
       [-0.49611623,  0.41854724, -2.01626424],
       [-1.53998828, -0.26894451,  0.39934594],
       [ 1.12051069, -0.05040955, -1.11466544],
       [-0.10370729, -0.79678004, -1.05270397]])
```

                                    (for each set individually)

```
>>> df = pd.DataFrame(np.random.randn(10, 3), columns=['A','B','C'])
>>> df
          A         B         C
0 -0.137115 -0.951781  2.441000
1  0.688829  0.566701  0.328983
2  0.660575  1.345944  0.359981
3  0.037090 -1.149348 -1.803005
4  0.605768  0.048554  0.603675
5  0.971914 -0.438938  0.336875
6 -0.496116  0.418547 -2.016264
7 -1.539988 -0.268945  0.399346
8  1.120511 -0.050410 -1.114665
9 -0.103707 -0.796780 -1.052704
```

*__btagging_nn.py__*

changes performed using
load_btagging_data
in python module
*__btag_nn__*

+ separate event weight arrays
(from (eta,pT) reweighting)
for training and validation set

prepared pandas.DataFrames:

**HDF**

• preprocessed data
• index lists for training, testing
• default dictionary

separation into
3 pandas.DataFrames:
• training
• validation (split from training set)
• testing

**python**

# Training

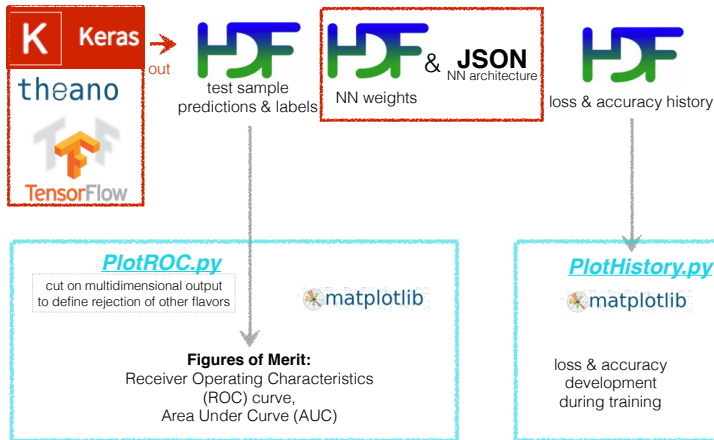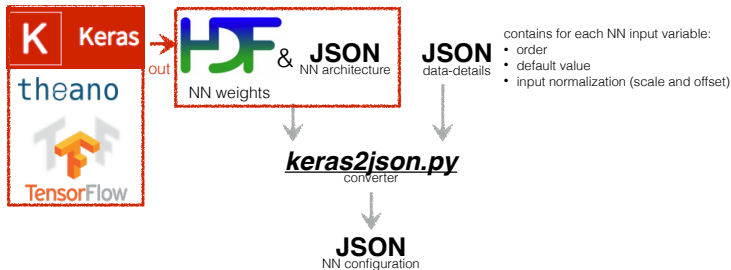# Plotting



**Select final training**

# Making the NN accessible in any C++ framework



**C++ lightweight classes to apply trained NN
in any C++ framework**
https://github.com/dguest/lwtnn.git

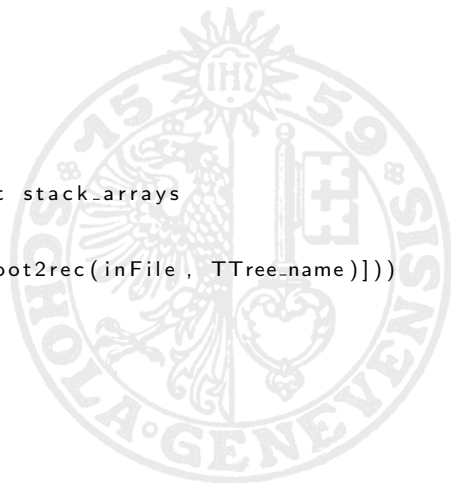(mainly developed by Dan Guest, UC Irvine)

# Converting flat ROOT ntuple into pandas.DataFrame

```python
import pandas as pd
from numpy.lib.recfunctions import stack_arrays
from root_numpy import root2rec

df = pd.DataFrame(stack_arrays([root2rec(inFile, TTree_name)]))
```
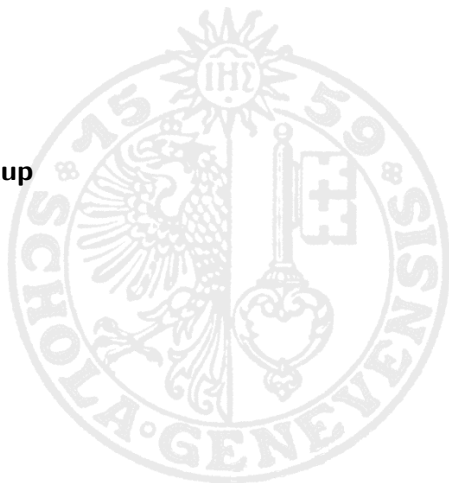
https://github.com/Marie89/BTagging_DL1.git

**Backup**

# How we do it in Athena (ATLAS C++ framework)