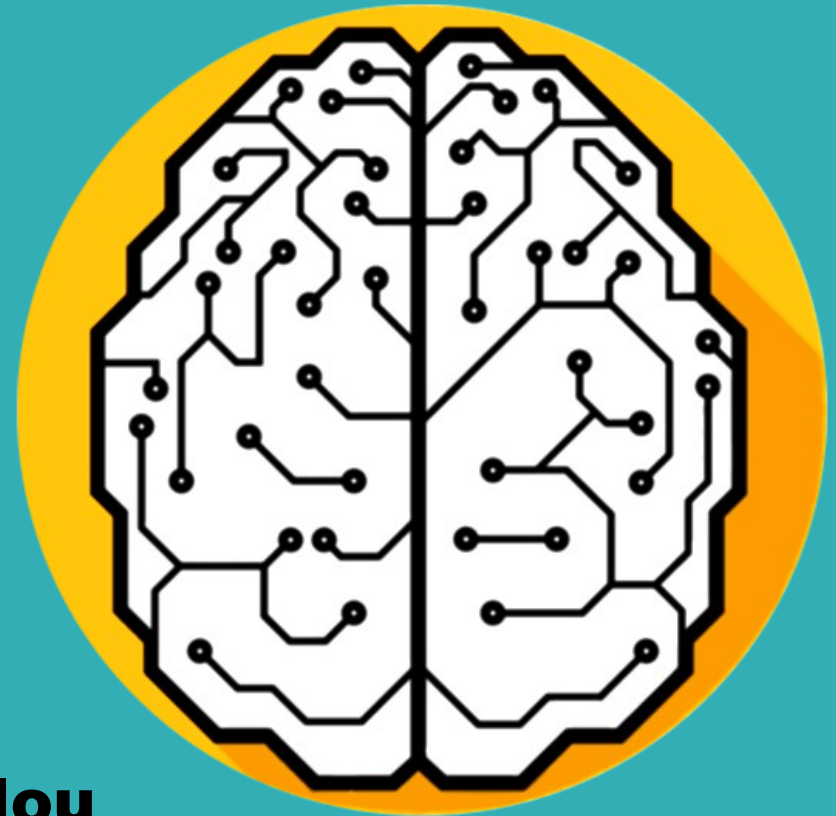


Parallelization of ROOT Machine Learning Methods



Pourya Vakilipourtakalou

Supervisors : Prof. Lorenzo Moneta

Prof. Sergei Gleyzer

Overview

- **Machine Learning**
- **ROOT**
- **TMVA**
- **Cross Validation**
- **Parallelization**
- **Outlook**



Machine Learning

Teaching the computers to do something exactly like the way people learn.



How do people learn?

H
o
r

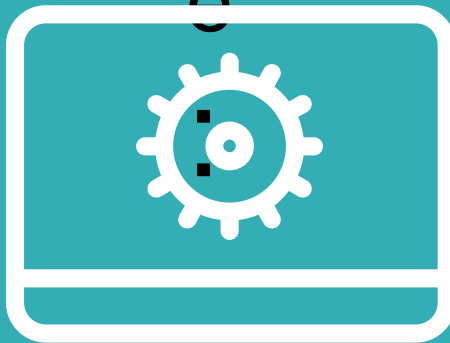


e
:



And this is Machine Learning!

C
a
k
e



Data
Algorithm



Output

Data
Output



Trained
Algorithm

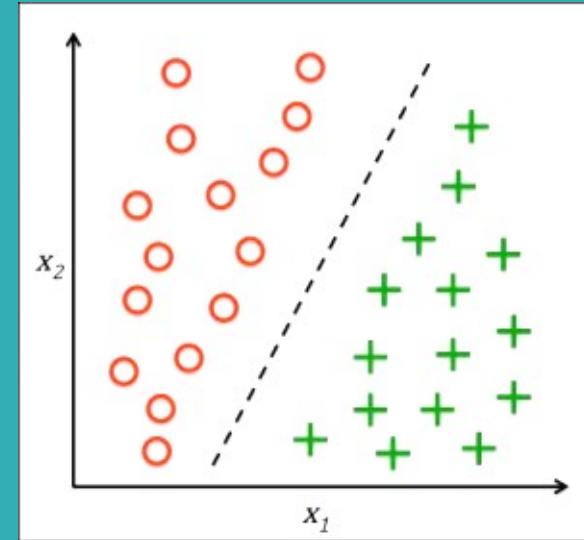


We
tra



Machine Learning

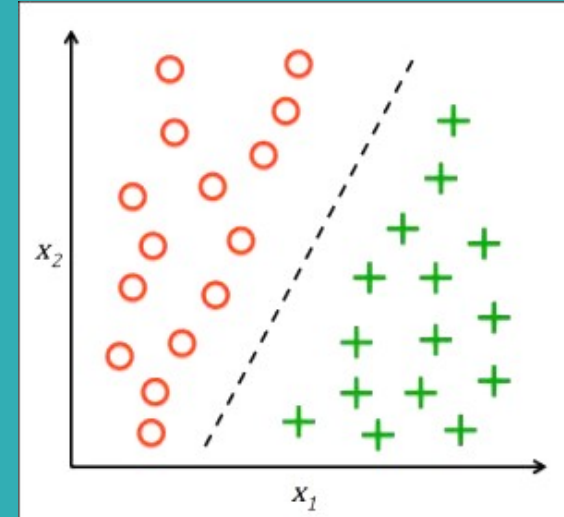
- An example of classification
- X_1 → age of the patient
- X_2 → size of the tumor
- Y → output : Malignant or Benign → 0 or 1
- Proposing a function like $H(X_1, X_2)$ like $aX_1 + bX_2$ → it can be anything
- Try to find optimal a and b



Machine Learning

More Physical Example

- $X \rightarrow$ vector of Kinematic Variables
- $Y \rightarrow$ output : Higgs (Signal) or Background $\rightarrow 0$ or 1
- Proposing a function like $F(X)$



m
a
i
n
l
y
w
r

m
o
d
u
l
a
r
s
c
i
e
n
t
i
c
R
O
O
T

f
u
n
c
t
i
o
n
a
l
i
t
y



ROOT
Data Analysis Framework



ROOT, Machine Learning → TMVA

- Toolkit for Multivariate Data Analysis
- Bunch of methods that provides a ROOT-integrated machine learning environment
- It includes Rectangular cut optimization, Boosted/Bagged decision trees, Artificial neural networks, ...



TMVA :

Training Phase → Training the method on known dataset

Application Phase → Applying the trained method to unknown dataset

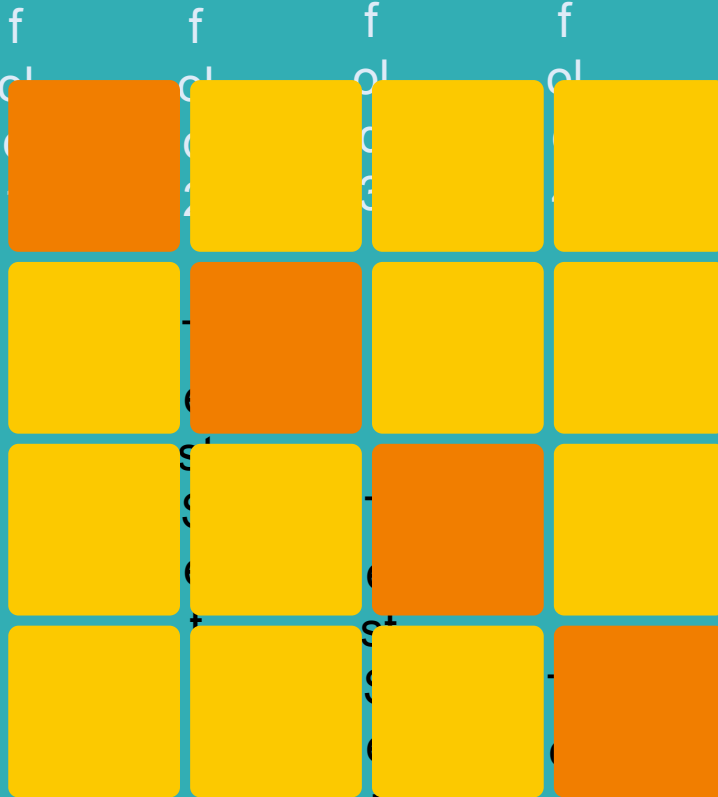


Training Phase

- Create the Factory → The connection between the user and TMVA
- Giving the Training/Test trees
- Register input variables (Features)
- Select the MVA methods from the Factory that we are going to train on the data set
- Book, Train and Test the Method



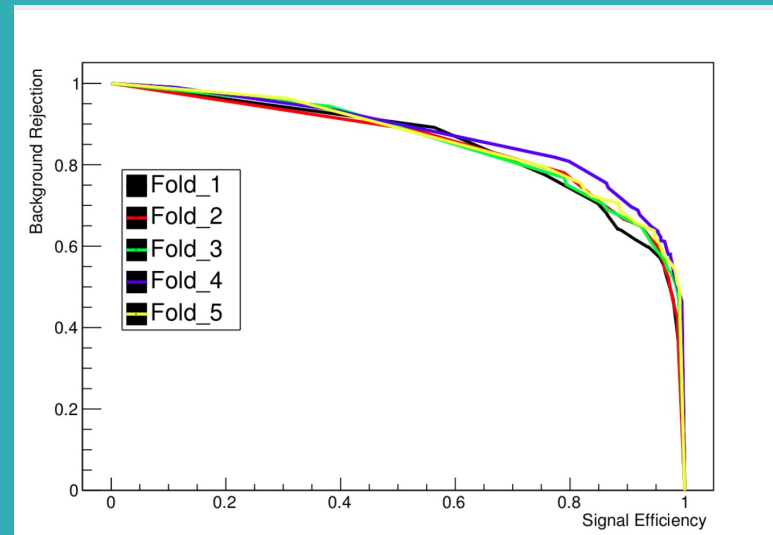
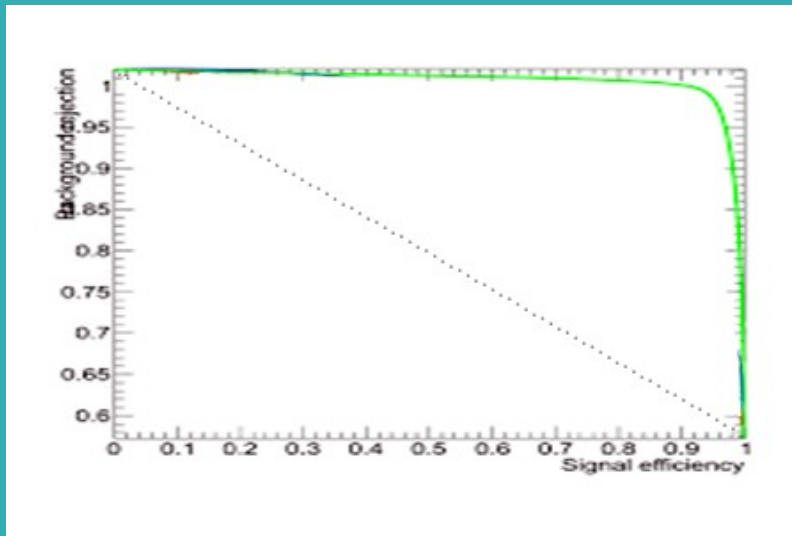
Cross Validation



- **T**raini**n**g d**a**t**a**set
 - **C**omple**t**e d**a**t**a**set
 - **T**est d**a**t**a**set
- C
a
l
c
u
l
a
t
e



Cross Validation: PlotROC()



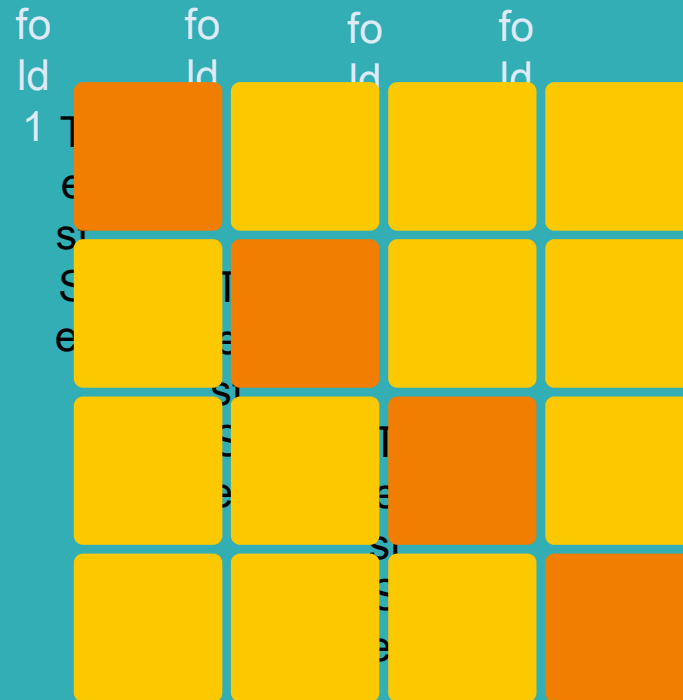
ROC → Receiver Operating Characteristic:
in Statistic graphical plot that illustrates the
performance of a classifier



Parallelization

ROOT Classes for Parallelization

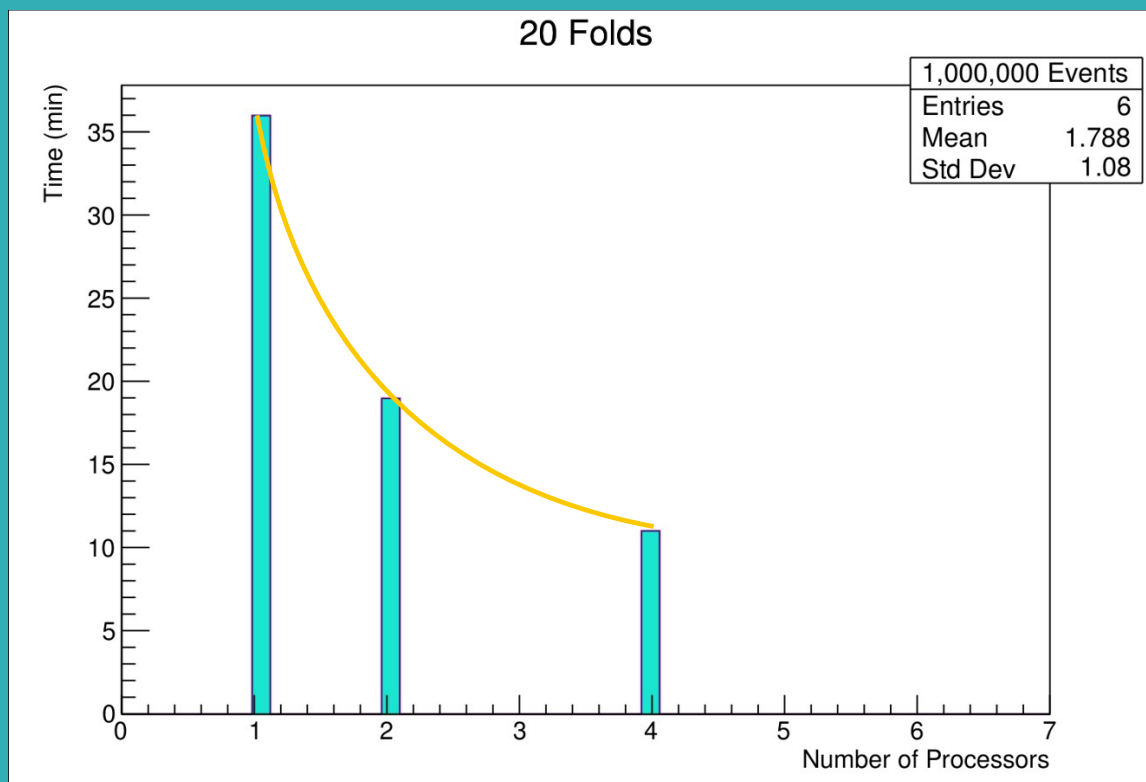
- **ThreadPool**
→ **Multithreading**
- **TProcPool**
→ **Multiprocessing**
- **Multithreading**
→ **More difficult to implement : needs locking → no Global**



This
S15

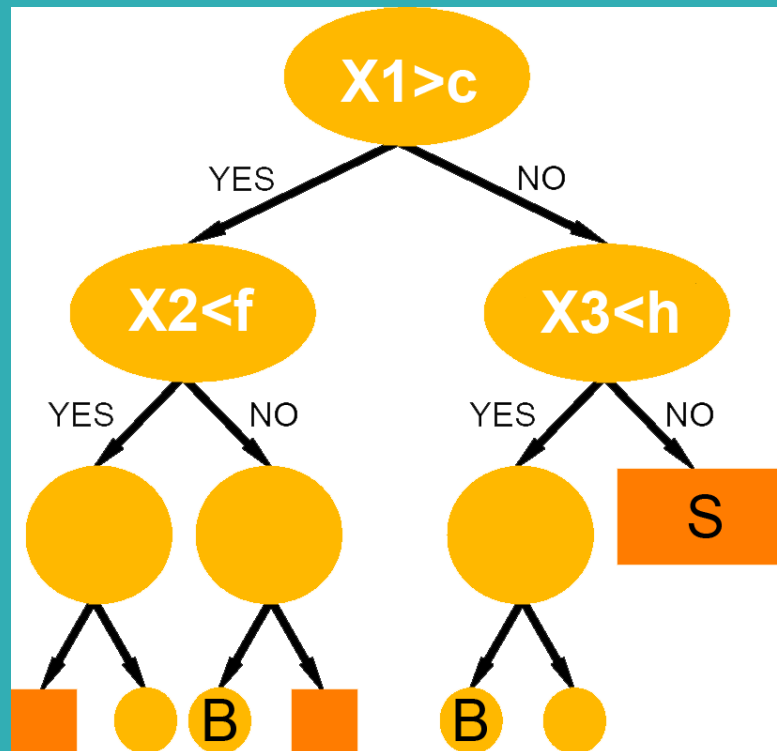
Thread 1
Thread 2
Thread 3
Thread 4





Outlook

Parallelization of different methods like
BDT → Boosted Decision Tree



Conclusion





Thank you all very
much for your
attention!

