

# Searches for the Higgs boson in the tau-tau decay channel

using CMS Open Data  
Final Project for CMEPDA exam

05/05/22

Domenico Riccardi<sup>1</sup> & Viola Floris<sup>2</sup>

[1] [domenico.riccardi@cern.ch](mailto:domenico.riccardi@cern.ch)

[2] [v.floris1@studenti.unipi.it](mailto:v.floris1@studenti.unipi.it)

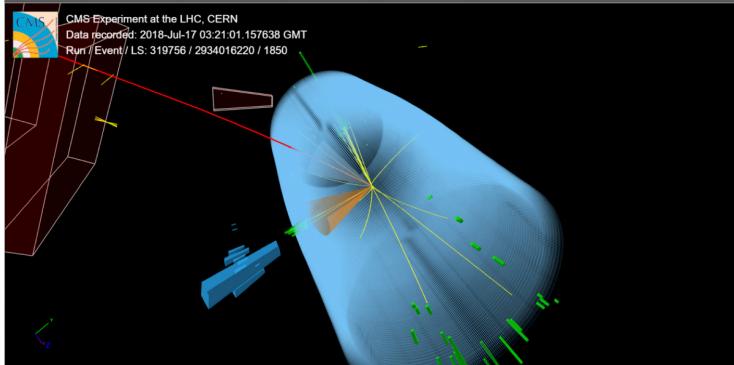
# GitHub Repository

<https://github.com/DomRic98/Searches-for-the-Higgs-boson-in-the-tau-tau-decay-channel>

README.md

build passing license GPL-3.0 languages 4 last commit today total lines 20k python 65.8%

## Searches for the Higgs boson in the tau-tau decay channel



CMS Experiment at the LHC, CERN  
Data recorded: 2018-Jul-17 03:21:01,157638 GMT  
Run / Event / LS: 319756 / 2934016220 / 1850

Event display of a p-p collision recorded at CMS with a candidate Higgs boson decaying into two tau leptons.  
[<https://cds.cern.ch/record/2725256>]

### Description

Welcome to this repository!

Here, you find all the material produced for the final project for Computing Methods for Experimental Physics and Data Analysis exam at the University of Pisa. Running the Python and C++ scripts that you can find in the folders, you can explore a little part of a more complex analysis on the search of Higgs boson decay into two tau leptons.

Environments 1

github-pages Active

Languages

Language	Percentage
Python	65.8%
C++	32.1%
Batchfile	1.2%
Makefile	0.9%

# Struttura del progetto

## Contents

### 1. Studio dei Dataset e selezione degli eventi (ROOT/C++)



- Fisica del processo
- Utilizzo di RDataFrame per lo sviluppo del codice di selezione
- Plotting delle distribuzioni delle variabili fisiche

### 2. Implementazione di algoritmi ML per la classificazione sig/bkg (Python)



- Artificial Neural Network
- Random Forest Classifier
- Plotting delle predizioni su Dataset di test
- Studio sulla selezione degli iperparametri delle reti
- Implementazione di vari unittest e documentazione del codice

### 3. Sviluppo di algoritmi per il download dei Dataset in sequenza e in parallelo



- Algoritmi e unittest

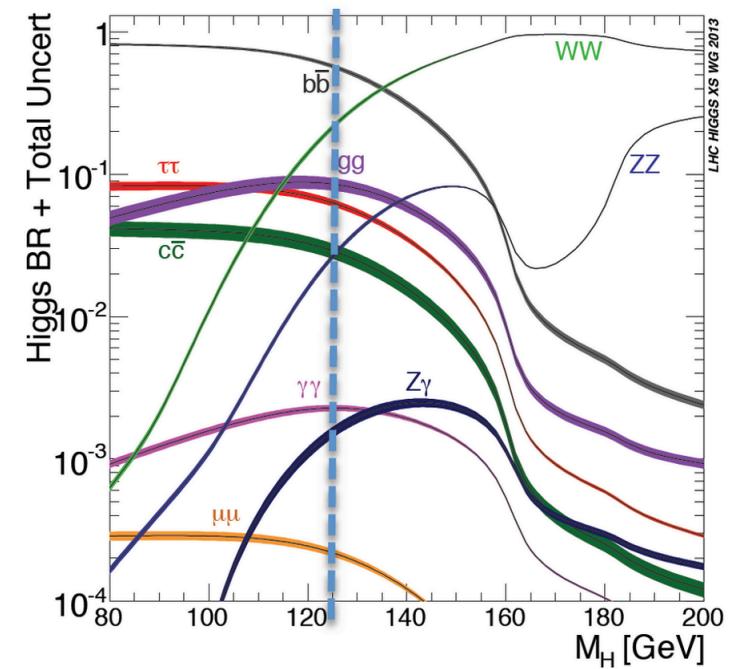
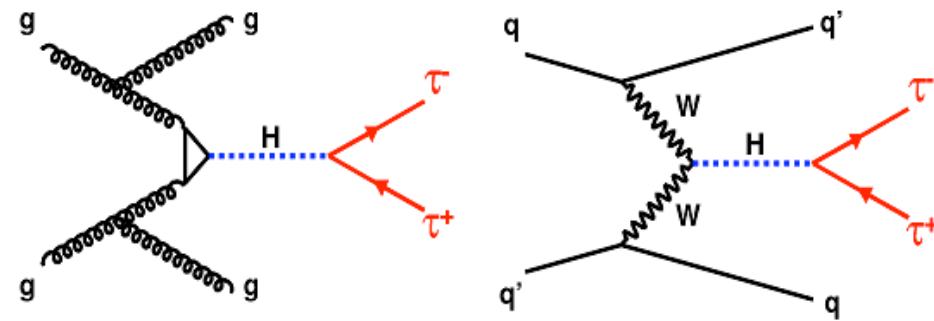
# **Studio dei Dataset e selezione degli eventi (ROOT/C++)**

Prima Parte



# Fisica del processo

- Obiettivo: Rivelare il decadimento  $H \rightarrow \tau_\mu \tau_h$ , dove con  $\tau_\mu$  si intende il processo  $\tau \rightarrow \mu \nu_\mu \nu_\tau$  e con  $\tau_h$  il conseguente decadimento in adroni
- Produzione dominante dell'Higgs a LHC:
  - Gluon gluon Fusion,  $\sigma(ggH) \sim 49 \text{ pb} @ 13\text{TeV}$
  - Vector Boson Fusion,  $\sigma(VBF) \sim 3.8 \text{ pb} @ 13\text{TeV}$

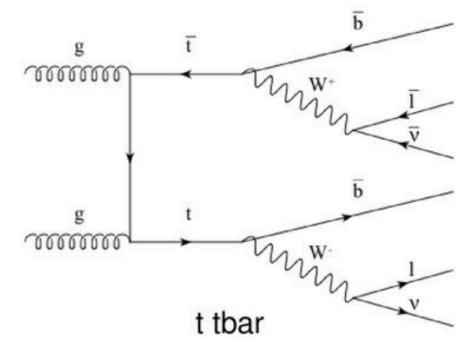
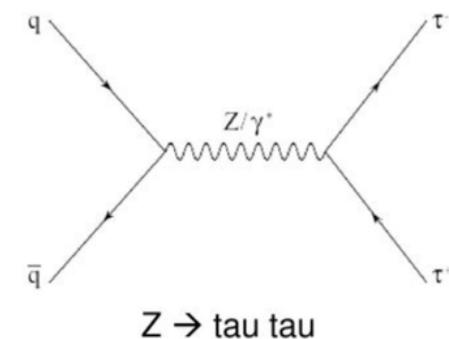


# Background e Dataset

- I fondi principali per il nostro canale di segnale sono:
  - Drell-Yan:  $Z/\gamma^* \rightarrow l^+l^-$
  - Produzione di coppie  $t\bar{t}$
  - $W + \text{jets}$

MC sample (ROOT files)	Cross section [pb]	Number of events
GluGluToHToTauTau	19.6	476963
VBF_HToTauTau	1.55	491653
DYJetsToLL	3503.7	30458871
TTbar	225.2	6423106
W1JetsToLNu	6381.2	29784800
W2JetsToLNu	2039.8	30693853
W3JetsToLNu	612.5	15241144

MC samples (CMS OpenData, 2012) @ 8 TeV  
con **sezione d'urto** e **numero di eventi** generati per canale



# CMS OpenData Dataset semantics

N	Variable	Type	Description
1	<b>event</b>	<i>unsigned long</i>	Event number
2	<b>HLT_IsoMu17_eta2p1_LooseIsolationPFTau20</b>	<i>bool</i>	Flag indicating whether the event passed the respective trigger
3	<b>PV_npv</b>	<i>int</i>	Number of primary vertices
4	<b>nMuon</b>	<i>unsigned int</i>	Number of muons in this event
5	<b>Muon_pt</b>	<i>float[nMuon]</i>	Transverse momentum of the muons
6	<b>Muon_eta</b>	<i>float[nMuon]</i>	Pseudorapidity of the muons
7	<b>Muon_phi</b>	<i>float[nMuon]</i>	Azimuth angle of the muons
8	<b>Muon_mass</b>	<i>float[nMuon]</i>	Mass of the muons
9	<b>Muon_charge</b>	<i>int[nMuon]</i>	Charge of the muons (either 1 or -1)
10	<b>Muon_tightId</b>	<i>bool[nMuon]</i>	Muon object identification with tight requirements
11	<b>nTau</b>	<i>unsigned int</i>	Number of taus in this event
12	<b>Tau_pt</b>	<i>float[nTau]</i>	Transverse momentum of the taus
13	<b>Tau_eta</b>	<i>float[nTau]</i>	Pseudorapidity of the taus
14	<b>Tau_phi</b>	<i>float[nTau]</i>	Azimuth angle of the taus

15	<b>Tau_mass</b>	<i>float[nTau]</i>	Mass of the taus
16	<b>Tau_charge</b>	<i>int[nTau]</i>	Charge of the taus (either 0, 1 or -1). 0 indicates that the charge could not be determined.
17	<b>Tau_decayMode</b>	<i>int[nTau]</i>	Hadronic decay mode of the taus
18	<b>Tau_reliso_all</b>	<i>float[nTau]</i>	Tau isolation divided by the transverse momentum
19	<b>Tau_idDecayMode</b>	<i>bool[nTau]</i>	Indicates passed decay mode finding of the taus
20	<b>Tau_idIso[VLoose, Loose, Medium, Tight]</b>	<i>bool[nTau]</i>	Indicates passed tau discriminator with [very loose, loose, medium, tight] requirements
21	<b>Tau_idAntiEle[Loose, Medium, Tight]</b>	<i>bool[nTau]</i>	Indicates passed anti-electron discriminator with [loose, medium, tight] requirements
22	<b>Tau_idAntiMu[Loose, Medium, Tight]</b>	<i>bool[nTau]</i>	Indicates passed anti-muon discriminator with [loose, medium, tight] requirements
23	<b>nJet</b>	<i>unsigned int</i>	Number of jets in this event
24	<b>Jet_pt</b>	<i>float[nJet]</i>	Transverse momentum of the jets
25	<b>Jet_eta</b>	<i>float[nJet]</i>	Pseudorapidity of the jets
26	<b>Jet_phi</b>	<i>float[nJet]</i>	Azimuth angle of the jets
27	<b>Jet_mass</b>	<i>float[nJet]</i>	Mass of the jets
28	<b>Jet_puld</b>	<i>bool[nJet]</i>	Jet pile-up identification flag
29	<b>Jet_btag</b>	<i>float[nJet]</i>	b-Jet identification score
30	<b>MET_pt</b>	<i>float</i>	Missing transverse energy
31	<b>MET_phi</b>	<i>float</i>	Azimuth angle of the missing transverse energy

GluGluToHtoTauTau

# Analysis code

- Abilita il multi-threading implicito di ROOT

```
int main()
{
    // Implicit parallelism in ROOT

    ROOT::EnableImplicitMT();
    int poolSize = ROOT::GetImplicitMPoolSize();
    std::cout << "Pool size used for multi-threading: "
        << poolSize << std::endl;

    std::vector<std::string> ROOT_files = {
        "GluGluToHToTauTau", "VBF_HToTauTau", // Signal MC
        "DYJetsToLL", "TTbar",
        "W1JetsToLNu", "W2JetsToLNu", "W3JetsToLNu", // Bkg MC
        "Run2012B_TauPlusX", "Run2012C_TauPlusX" // Real Data
    };

    auto start = std::chrono::system_clock::now();

    for (const std::string &sample : ROOT_files) {

        std::cout << "*****" << std::endl;
        std::cout << "Processing sample: " << sample << std::endl;
        ROOT::RDataFrame d_root("Events", sample + ".root");

        int num_events = *d_root.Count();
        const float integratedLuminosity = 11.467 * 1000.0; // (pb^-1)

        float weight = 1.;
        if(sample.find("Run") != 0)
            weight = (cross_section[sample]/num_events) * integratedLuminosity;
        auto d_MCweight = d_root.Define("weight", [weight](){ return weight;});

        std::cout << "\t Number of events in ROOT file: " << num_events
            << std::endl;

        auto d_muon = muon_selection(d_MCweight);
        std::cout << "\t Number of events after muon selection: "
            << *d_muon.Count() << std::endl;

        auto d_tau = tau_selection(d_muon);
        std::cout << "\t Number of events after tau selection: "
            << *d_tau.Count() << std::endl;

        auto d_h2mutau = h2mutau(d_tau);
        std::cout << "\t Number of events after selecting candidates for the
            Higgs (H->tau): " << *d_h2mutau.Count() << std::endl;

        auto d_variables = ML_variables(d_h2mutau);
        auto d_jets = take_jet(d_variables);
        auto d = d_jets.Filter("muon_charge*tau_charge<0");

        // Save the new DataFrame of the events selected in the analysis
        // procedure.
        d.Snapshot("Events", sample + "_selected.root", finalVariables);
    }

    auto end = std::chrono::system_clock::now();
    std::chrono::duration<double> elapsed_seconds = end-start;
    std::time_t end_time = std::chrono::system_clock::to_time_t(end);
    std::cout << "Finished computation at "
        << std::ctime(&end_time) << "Total elapsed time: "
        << elapsed_seconds.count() << "s" << std::endl;
}
```

# Analysis code

- Definisco una nuova colonna del DataFrame: peso dell'evento MC

```
int main()
{
    // Implicit parallelism in ROOT

    ROOT::EnableImplicitMT();
    int poolSize = ROOT::GetImplicitMPoolSize();
    std::cout << "Pool size used for multi-threading: "
        << poolSize << std::endl;

    std::vector<std::string> ROOT_files = {
        "GluGluToHToTauTau", "VBF_HToTauTau", // Signal MC
        "DYJetsToLL", "TTbar",
        "W1JetsToLNu", "W2JetsToLNu", "W3JetsToLNu", // Bkg MC
        "Run2012B_TauPlusX", "Run2012C_TauPlusX" // Real Data
    };

    auto start = std::chrono::system_clock::now();

    for (const std::string &sample : ROOT_files) {
        std::cout << "*****" << std::endl;
        std::cout << "Processing sample: " << sample << std::endl;
        ROOT::RDataFrame d_root("Events", sample + ".root");

        int num_events = *d_root.Count();
        const float integratedLuminosity = 11.467 * 1000.0; // (pb^-1)

        float weight = 1.;
        if(sample.find("Run") != 0)
            weight = (cross_section[sample]/num_events) * integratedLuminosity;
        auto d_MCweight = d_root.Define("weight", [weight](){ return weight;});

        std::cout << "\t Number of events in ROOT file: " << num_events
            << std::endl;

        auto d_muon = muon_selection(d_MCweight);
        std::cout << "\t Number of events after muon selection: "
            << *d_muon.Count() << std::endl;

        auto d_tau = tau_selection(d_muon);
        std::cout << "\t Number of events after tau selection: "
            << *d_tau.Count() << std::endl;

        auto d_h2mutau = h2mutau(d_tau);
        std::cout << "\t Number of events after selecting candidates for the
            Higgs (H->tau): " << *d_h2mutau.Count() << std::endl;

        auto d_variables = ML_variables(d_h2mutau);
        auto d_jets = take_jet(d_variables);
        auto d = d_jets.Filter("muon_charge*tau_charge<0");

        // Save the new DataFrame of the events selected in the analysis
        // procedure.
        d.Snapshot("Events", sample + "_selected.root", finalVariables);
    }

    auto end = std::chrono::system_clock::now();
    std::chrono::duration<double> elapsed_seconds = end-start;
    std::time_t end_time = std::chrono::system_clock::to_time_t(end);
    std::cout << "Finished computation at "
        << std::ctime(&end_time) << "Total elapsed time: "
        << elapsed_seconds.count() << "s" << std::endl;
}
```

# Analysis code

- Flusso delle chiamate alle funzioni di selezione & salvataggio del nuovo RDataFrame come ROOT file

```
int main()
{
    // Implicit parallelism in ROOT

    ROOT::EnableImplicitMT();
    int poolSize = ROOT::GetImplicitMPoolSize();
    std::cout << "Pool size used for multi-threading: "
        << poolSize << std::endl;

    std::vector<std::string> ROOT_files = {
        "GluGluToHToTauTau", "VBF_HToTauTau", // Signal MC
        "DYJetsToLL", "TTbar",
        "W1JetsToLNu", "W2JetsToLNu", "W3JetsToLNu", // Bkg MC
        "Run2012B_TauPlusX", "Run2012C_TauPlusX" // Real Data
    };

    auto start = std::chrono::system_clock::now();

    for (const std::string &sample : ROOT_files) {

        std::cout << "*****" << std::endl;
        std::cout << "Processing sample: " << sample << std::endl;
        ROOT::RDataFrame d_root("Events", sample + ".root");

        int num_events = *d_root.Count();
        const float integratedLuminosity = 11.467 * 1000.0; // (pb^-1)

        float weight = 1.;
        if(sample.find("Run") != 0)
            weight = (cross_section[sample]/num_events) * integratedLuminosity;
        auto d_MCweight = d_root.Define("weight", [weight](){ return weight;});
    }
}
```

```
    std::cout << "\t Number of events in ROOT file: " << num_events
        << std::endl;

    auto d_muon = muon_selection(d_MCweight);
    std::cout << "\t Number of events after muon selection: "
        << *d_muon.Count() << std::endl;

    auto d_tau = tau_selection(d_muon);
    std::cout << "\t Number of events after tau selection: "
        << *d_tau.Count() << std::endl;

    auto d_h2mutau = h2mutau(d_tau);
    std::cout << "\t Number of events after selecting candidates for the
        Higgs (H->tau): " << *d_h2mutau.Count() << std::endl;

    auto d_variables = ML_variables(d_h2mutau);
    auto d_jets = take_jet(d_variables);
    auto d = d_jets.Filter("muon_charge*tau_charge<0");

    // Save the new DataFrame of the events selected in the analysis
    // procedure.
    d.Snapshot("Events", sample + "_selected.root", finalVariables);

}

auto end = std::chrono::system_clock::now();
std::chrono::duration<double> elapsed_seconds = end-start;
std::time_t end_time = std::chrono::system_clock::to_time_t(end);
std::cout << "Finished computation at "
    << std::ctime(&end_time) << "Total elapsed time: "
    << elapsed_seconds.count() << "s" << std::endl;
```

# Analysis code

## Esempio di output

- Eseguo dal prompt ROOT

```
domenico@MacBook-Pro-di-domenico ROOT_workspace % root -l
root [0] .L analysis.cpp
root [1] main()
```

- Oppure produco e run l'eseguibile con

```
domenico@MacBook-Pro-di-domenico ROOT_workspace % g++ analysis.cpp `root-config --cflags --glibs` -o analysis
domenico@MBP-di-domenico ROOT_workspace % ./analysis
```

```
Pool size used for multi-threading: 4
*****
Processing sample: GluGluToHToTauTau
Number of events in ROOT file: 476963
Number of events after muon selection: 32196
Number of events after tau selection: 8085
Number of events after selecting candidates for the Higgs (H->ττ): 8073
*****
Processing sample: VBF_HToTauTau
Number of events in ROOT file: 491653
Number of events after muon selection: 47233
Number of events after tau selection: 10493
Number of events after selecting candidates for the Higgs (H->ττ): 10441
*****
Processing sample: DYJetsToLL
Number of events in ROOT file: 30458871
Number of events after muon selection: 4704720
Number of events after tau selection: 70138
Number of events after selecting candidates for the Higgs (H->ττ): 66311
*****
Processing sample: TBar
Number of events in ROOT file: 6423106
Number of events after muon selection: 848658
Number of events after tau selection: 36921
Number of events after selecting candidates for the Higgs (H->ττ): 35324
*****
Processing sample: W1JetsToLNu
Number of events in ROOT file: 29784800
Number of events after muon selection: 1208333
Number of events after tau selection: 24015
Number of events after selecting candidates for the Higgs (H->ττ): 23962
*****
Processing sample: W2JetsToLNu
Number of events in ROOT file: 30693853
Number of events after muon selection: 2281445
Number of events after tau selection: 42476
Number of events after selecting candidates for the Higgs (H->ττ): 42341
*****
Processing sample: W3JetsToLNu
Number of events in ROOT file: 15241144
Number of events after muon selection: 1467529
Number of events after tau selection: 23989
Number of events after selecting candidates for the Higgs (H->ττ): 23889
*****
Processing sample: Run2012B_TauPlusX
Number of events in ROOT file: 35647508
Number of events after muon selection: 9028134
Number of events after tau selection: 136411
Number of events after selecting candidates for the Higgs (H->ττ): 133783
*****
Processing sample: Run2012C_TauPlusX
Number of events in ROOT file: 51303171
Number of events after muon selection: 12578182
Number of events after tau selection: 210897
Number of events after selecting candidates for the Higgs (H->ττ): 206812
Finished computation on Tue Apr 26 16:32:57 2022
Total elapsed time: 3418.47s
```

# Analysis code

## Selection functions

```
// namespace-alias-definition
using RFilter =
    ROOT::RDF::RInterface<ROOT::Detail::RDF::RJittedFilter>;

/*
* Perform a muons selection by selecting the final state with
* at least one good  $\mu$  candidate.
*/
template <typename T>
RFilter muon_selection(T &d)
{
    auto d_temp = d.Filter("nMuon>0 &
                           HLT_IsoMu17_eta2p1_LooseIsoPFTau20==1")
               .Define("goodMuons", "(abs(Muon_eta)<2.1) &
                        (Muon_pt>17) & (Muon_tightId==1)")
               .Filter("ROOT::VecOps::Sum(goodMuons)>0");
    return d_temp;
}

/*
 * Find at least one good tau.
 */
template <typename T>
RFilter tau_selection(T &d)
{
    auto d_temp = d.Filter("nTau>0 &
                           HLT_IsoMu17_eta2p1_LooseIsoPFTau20==1")
               .Define("goodTaus", "Tau_charge!=0 &
                           abs(Tau_eta)<2.3 & Tau_pt>20 &
                           Tau_idDecayMode==1 &
                           Tau_idIsoTight==1 &
                           Tau_idAntiEleTight==1 &
                           Tau_idAntiMuTight==1")
               .Filter("ROOT::VecOps::Sum(goodTaus)>0");
    return d_temp;
}
```

# Analysis code

- Definisco vari constraints sulle variabili del jet

```
/*
 * Select a "good jet" in the sample.
 */
template <typename T>
RFilter take_jet(T &d)
{
    auto first_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                          ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>0)
            return var[0];
        return -999.f;
    };

    auto second_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                          ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>1)
            return var[1];
        return -999.f;
    };

    auto compute_mjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.M());
        }
        return -999.f;
    };

    auto compute_ptjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.Pt());
        }
    };

    auto compute_jdata = [](int nGoodJets, float jeta_1, float jeta_2){
        if (nGoodJets>1) {
            return jeta_1 - jeta_2;
        }
        return -999.f;
    };

    auto d_temp = d.Define("goodJets", "Jet_pt>30 & abs(Jet_eta)<4.7 & Jet_puId==1")
        .Define("goodJets_idx", "ROOT::VecOps::Nonzero(goodJets)")
        .Define("jpt_1", first_jet, {"goodJets_idx", "Jet_pt"})
        .Define("jpt_2", second_jet, {"goodJets_idx", "Jet_pt"})
        .Define("jeta_1", first_jet, {"goodJets_idx", "Jet_eta"})
        .Define("jeta_2", second_jet, {"goodJets_idx", "Jet_eta"})
        .Define("jphi_1", first_jet, {"goodJets_idx", "Jet_phi"})
        .Define("jphi_2", second_jet, {"goodJets_idx", "Jet_phi"})
        .Define("jm_1", first_jet, {"goodJets_idx", "Jet_mass"})
        .Define("jm_2", second_jet, {"goodJets_idx", "Jet_mass"})
        .Define("jbtags_1", first_jet, {"goodJets_idx", "Jet_btag"})
        .Define("jbtags_2", second_jet, {"goodJets_idx", "Jet_btag"})
        .Define("nGoodJets", "int(goodJets_idx.size())")
        .Define("jp4_1", "ROOT::Math::PtEtaPhiMVector(jpt_1, jeta_1,
                                                    jphi_1, jm_1)")
        .Define("jp4_2", "ROOT::Math::PtEtaPhiMVector(jpt_2, jeta_2,
                                                    jphi_2, jm_2)")
        .Define("jjp4", "jp4_1+jp4_2")
        .Define("jj_m", compute_mjj, {"nGoodJets", "jjp4"})
        .Define("jj_pt", compute_ptjj, {"nGoodJets", "jjp4"})
        .Define("jj_delta", compute_jdata, {"nGoodJets", "jeta_1", "jeta_2"});
    return d_temp;
}
```

# Analysis code

- Nuove colonne nel DataFrame per il primo e il secondo jet

```
/*
 * Select a "good jet" in the sample.
 */
template <typename T>
RFilter take_jet(T &d)
{
    auto first_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                          ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>0)
            return var[0];
        return -999.f;
    };

    auto second_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                           ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>1)
            return var[1];
        return -999.f;
    };

    auto compute_mjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.M());
        }
        return -999.f;
    };

    auto compute_ptjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.Pt());
        }
    };

    auto jeta_1 = [] (int nGoodJets, float jeta_1, float jeta_2)
    {
        if (nGoodJets>1) {
            return jeta_1 - jeta_2;
        }
        return -999.f;
    };

    auto jeta_2 = [] (int nGoodJets, float jeta_1, float jeta_2)
    {
        if (nGoodJets>1) {
            return jeta_2 - jeta_1;
        }
        return -999.f;
    };

    auto d_temp = d.Define("goodJets", "Jet_pt>30 & abs(Jet_eta)<4.7 & Jet_puId==1")
        .Define("goodJets_idx", "ROOT::VecOps::Nonzero(goodJets)")
        .Define("jpt_1", first_jet, {"goodJets_idx", "Jet_pt"})
        .Define("jpt_2", second_jet, {"goodJets_idx", "Jet_pt"})
        .Define("jeta_1", first_jet, {"goodJets_idx", "Jet_eta"})
        .Define("jeta_2", second_jet, {"goodJets_idx", "Jet_eta"})
        .Define("jphi_1", first_jet, {"goodJets_idx", "Jet_phi"})
        .Define("jphi_2", second_jet, {"goodJets_idx", "Jet_phi"})
        .Define("jm_1", first_jet, {"goodJets_idx", "Jet_mass"})
        .Define("jm_2", second_jet, {"goodJets_idx", "Jet_mass"})
        .Define("jbtags_1", first_jet, {"goodJets_idx", "Jet_btag"})
        .Define("jbtags_2", second_jet, {"goodJets_idx", "Jet_btag"})
        .Define("nGoodJets", "int(goodJets_idx.size())")
        .Define("jp4_1", "ROOT::Math::PtEtaPhiMVector(jpt_1, jeta_1,
                                                     jphi_1, jm_1)")
        .Define("jp4_2", "ROOT::Math::PtEtaPhiMVector(jpt_2, jeta_2,
                                                     jphi_2, jm_2)")
        .Define("jjp4", "jp4_1+jp4_2")
        .Define("jj_m", compute_mjj, {"nGoodJets", "jjp4"})
        .Define("jj_pt", compute_ptjj, {"nGoodJets", "jjp4"})
        .Define("jj_delta", compute_jdata, {"nGoodJets", "jeta_1", "jeta_2"});
    return d_temp;
}
```

# Analysis code

- Creo oggetti di tipo 4-vettore con PtEtaPhiMVector() e calcolo nuove variabili

```
/*
 * Select a "good jet" in the sample.
 */
template <typename T>
RFilter take_jet(T &d)
{
    auto first_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                          ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>0)
            return var[0];
        return -999.f;
    };

    auto second_jet = [] (ROOT::RVec<unsigned long> &goodJets_idx,
                          ROOT::RVec<float> &jet_var)
    {
        auto var = ROOT::VecOps::Take(jet_var, goodJets_idx);
        if (goodJets_idx.size()>1)
            return var[1];
        return -999.f;
    };

    auto compute_mjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.M());
        }
        return -999.f;
    };

    auto compute_ptjj = [] (int nGoodJets, ROOT::Math::PtEtaPhiMVector jjp4)
    {
        if (nGoodJets>1) {
            return float(jjp4.Pt());
        }
    };
}
```

```
    }

    return -999.f;
};

auto compute_jdata = [] (int nGoodJets, float jeta_1, float jeta_2){
    if (nGoodJets>1) {
        return jeta_1 - jeta_2;
    }
    return -999.f;
};

auto d_temp = d.Define("goodJets", "Jet_pt>30 & abs(Jet_eta)<4.7 & Jet_puId==1")
    .Define("goodJets_idx", "ROOT::VecOps::Nonzero(goodJets)")
    .Define("jpt_1", first_jet, {"goodJets_idx", "Jet_pt"})
    .Define("jpt_2", second_jet, {"goodJets_idx", "Jet_pt"})
    .Define("jeta_1", first_jet, {"goodJets_idx", "Jet_eta"})
    .Define("jeta_2", second_jet, {"goodJets_idx", "Jet_eta"})
    .Define("jphi_1", first_jet, {"goodJets_idx", "Jet_phi"})
    .Define("jphi_2", second_jet, {"goodJets_idx", "Jet_phi"})
    .Define("jm_1", first_jet, {"goodJets_idx", "Jet_mass"})
    .Define("jm_2", second_jet, {"goodJets_idx", "Jet_mass"})
    .Define("jbtag_1", first_jet, {"goodJets_idx", "Jet_btag"})
    .Define("jbtag_2", second_jet, {"goodJets_idx", "Jet_btag"})
    .Define("nGoodJets", "int(goodJets_idx.size())")
    .Define("jp4_1", "ROOT::Math::PtEtaPhiMVector(jpt_1, jeta_1,
                                                jphi_1, jm_1)")
    .Define("jp4_2", "ROOT::Math::PtEtaPhiMVector(jpt_2, jeta_2,
                                                jphi_2, jm_2)")
    .Define("jjp4", "jp4_1+jp4_2")
    .Define("jj_m", compute_mjj, {"nGoodJets", "jjp4"})
    .Define("jj_pt", compute_ptjj, {"nGoodJets", "jjp4"})
    .Define("jj_delta", compute_jdata, {"nGoodJets", "jeta_1", "jeta_2"});
}

return d_temp;
}
```

# Analysis code

- Trovo la coppia che rappresenta un “buon” candidato Higgs:  $H \rightarrow \tau_\mu \tau_h$

```
/*
 * Define a function which constructs a muon-tau pair from the
 * collections of goodMuons and goodTaus.
 */
template <typename T>
RFilter h2mutau(T &d)
{
    auto goodPair = [] (ROOT::RVec<int> &goodMuons,
                        ROOT::RVec<int> &goodTaus,
                        ROOT::RVec<float> &eta1,
                        ROOT::RVec<float> &eta2,
                        ROOT::RVec<float> &phi1,
                        ROOT::RVec<float> &phi2,
                        ROOT::RVec<float> &pt1,
                        ROOT::RVec<float>& iso2,
                        ROOT::RVec<int>& charge1,
                        ROOT::RVec<int>& charge2)
    {
        auto muons_idx = ROOT::VecOps::Nonzero(goodMuons); //idx of μ
        auto taus_idx = ROOT::VecOps::Nonzero(goodTaus); //idx of τ
        // Find valid pairs based on delta r
        if(taus_idx.size() == 1 & muons_idx.size() == 1){
            // The function DeltaR computes ΔR=√((η1-η2)²+(φ1-φ2)² of
            // the given collections eta1, eta2, phi1 and phi2
            auto DeltaR = ROOT::VecOps::DeltaR(
                eta1[muons_idx[0]], eta2[taus_idx[0]],
                phi1[muons_idx[0]], phi2[taus_idx[0]]);
            // Select the pair with ΔR > 0.5
            if(DeltaR > 0.5){
                return std::pair<int, int>(muons_idx[0], taus_idx[0]);
            }
            else{
                // If there is more than one good muon or good tau, find
                // the muon with the maximum pt and the tau with the
                // minimum value of isolation.
                auto idx_pt_max = ROOT::VecOps::ArgMax(pt1);
                auto idx_iso_min = ROOT::VecOps::ArgMin(iso2);
                auto DeltaR = ROOT::VecOps::DeltaR(
                    eta1[idx_pt_max], eta2[idx_iso_min],
                    phi1[idx_pt_max], phi2[idx_iso_min]);
                if(DeltaR > 0.5){
                    return std::pair<int, int>(idx_pt_max, idx_iso_min);
                }
                return std::pair<int, int>(-999, -999);
            }
        };
        auto d_temp = d.Define("h2mutau", goodPair, {"goodMuons",
                                                       "goodTaus", "Muon_eta", "Tau_eta", "Muon_phi", "Tau_phi",
                                                       "Muon_pt", "Tau_relIso_all", "Muon_charge", "Tau_charge"})
                     .Define("idx_muon", "h2mutau.first")
                     .Filter("idx_muon >= 0")
                     .Define("idx_tau", "h2mutau.second");
        return d_temp;
    }
}
```

# Analysis code

- Definisco (e calcolo) le variabili del nuovo RDataFrame che utilizzeremo per il Machine Learning

```
template <typename T>
RFilter ML_variables(T &d)
{
    // Lambda function to calculate the transverse mass between
    // lepton and MET.
    auto compute_mt = [](float pt, float phi,
                         float met_pt, float met_phi){
        double dphi = ROOT::VecOps::DeltaPhi(phi, met_phi);
        return std::sqrt(2.0 * pt * met_pt * (1.0 -
            std::cos(dphi)));
    };

    auto compute_DeltaR = [](float mu_eta, float tau_eta,
                           float mu_phi, float tau_phi){
        return ROOT::VecOps::DeltaR(mu_eta, tau_eta,
                                   mu_phi, tau_phi);
    };

    auto d_temp = d.Define("muon_pt", "Muon_pt[idx_muon]")
                  .Define("muon_eta", "Muon_eta[idx_muon]")
                  .Define("muon_phi", "Muon_phi[idx_muon]")
                  .Define("muon_m", "Muon_mass[idx_muon]")
                  .Define("muon_charge", "Muon_charge[idx_muon]")
                  .Define("muon_p4",
                         "ROOT::Math::PtEtaPhiMVector(muon_pt,
                         muon_eta, muon_phi, muon_m)")
                  .Define("muon_iso",
                         "Muon_pfRelIso03_all[idx_muon]")
                  .Define("tau_pt", "Tau_pt[idx_tau]")
                  .Define("tau_eta", "Tau_eta[idx_tau]")
                  .Define("tau_phi", "Tau_phi[idx_tau]")
                  .Define("tau_m", "Tau_mass[idx_tau]")
                  .Define("tau_charge", "Tau_charge[idx_muon]")
                  .Define("tau_iso", "Tau_relIso_all[idx_tau]")
                  .Define("tau_p4", "ROOT::Math::PtEtaPhiMVector
                         (tau_pt, tau_eta, tau_phi, tau_m)")
                  .Define("p4", "muon_p4+tau_p4")
                  .Define("m_vis", "p4.M()")
                  .Define("pt_vis", "p4.Pt()")
                  .Define("mt_mu", compute_mt,
                         {"muon_pt", "muon_phi", "MET_pt", "MET_phi"})
                  .Define("mt_tau", compute_mt,
                         {"tau_pt", "tau_phi", "MET_pt", "MET_phi"})
                  .Define("dRmu_tau", compute_DeltaR,
                         {"muon_eta", "tau_eta", "muon_phi", "tau_phi"});
    }
    return d_temp;
}
```

# Plotting code

- Main function per la produzione delle distribuzioni delle variabili elencate prima

```
int main()
{
    // Create a new folder to contain the distribution plots produced
    std::string dirName = "Distribution_plots";
    std::stringstream bufH;
    bufH << dirName;
    if (!std::experimental::filesystem::exists(bufH.str()))
    {
        std::experimental::filesystem::create_directories(bufH.str());
    }

    // Define a std::map that contains the pair sample-RDataFrame and a
    // std::vector of ROOT files' names
    std::map<std::string, RFilter> samples;
    std::vector<std::string> ROOT_files = {
        "GluGluToHToTauTau_selected",
        "VBF_HToTauTau_selected",
        "DYJetsToLL_selected",
        "TTbar_selected",
        "W*JetsToLNu_selected",
        "Run2012*_TauPlusX_selected"
    };

    // For loop to process the events in ROOT files and to apply final
    // cuts to reduce the background (mostly W+jets)
    for (const std::string &channel: ROOT_files) {
        std::cout << "\tProcessing: " << channel << std::endl;
        ROOT::RDataFrame d("Events", channel + ".root");
    }
}
```

```
auto d_filter = d.Filter("mt_mu<30 & muon_iso<0.1")
    .Filter("tau_pt>25 & muon_pt>20");
samples.insert(std::pair<std::string,RFilter>(channel,d_filter));
std::cout << "\tNumber of events: "
    << *samples.at(channel).Count() << std::endl;
}

// Plot distribution of physical variables by using the plotting
// function
std::cout << " Plotting distribution" << std::endl;
plotting(samples, "m_vis", "Visible mass m_{#mu}+m_{#tau}",
    "m_{vis}", "[GeV]", 30, 20, 140);
plotting(samples, "pt_vis", "Visible transverse momentum",
    "p_{T,vis}", "[GeV]", 30, 0, 60);
plotting(samples, "tau_pt", "#tau transverse momentum", "p_{T
    (#tau )}", "[GeV]", 30, 17, 70);
plotting(samples, "muon_pt", "#mu transverse momentum", "p_{T} (#mu)",
    "[GeV]", 30, 17, 70);
plotting(samples, "muon_eta", "#mu pseudorapidity", "#eta(#mu)",
    "", 30, -2.1, 2.1);

// And other 28 variables!
}
```

# Plotting code

- Creo istogrammi per ciascun sample

```
template <typename T>
void plotting(T samples, std::string variable, const char *title, std::string xlabel_,
std::string unit, int bins, float lim_left, float lim_right) {
    std::cout << "\tPlotting: " << variable << std::endl;

    TCanvas *c = new TCanvas("c","c", 1000, 700);
    gStyle->SetOptStat(0);
    gStyle->SetTitleFontSize(0.05);
    // Define six 1D histograms for each channel
    auto h1 = samples.at("GluGluToHToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h2 = samples.at("VBF_HToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h3 = samples.at("DYJetsToLL_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h4 = samples.at("TTbar_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h5 = samples.at("WJetsToLNu_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h6 = samples.at("Run2012*TauPlusX_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable);

    std::vector<decltype(h1)> hist = {h1, h2, h3, h4, h5};
    std::vector<decltype(kRed+1)> color =
        {kRed+1, kMagenta-9, kGreen-3, kAzure+7, kOrange-3};
    // Assign the colours to each histogram
    for (int i=0; i<hist.size(); i++) {
        hist[i]->SetLineColor(color[i]);
        hist[i]->SetFillColor(color[i]);
    }
    // Create a stacked histogram with MC simulation samples
    THStack *hs = new THStack("hs", "");
    for (auto &h : hist) {
        hs->Add(h.GetPtr());
    }
    hs->Draw("HIST0");
    // Draw the Data with black points and error bars
    h6->SetMarkerStyle(kFullCircle);
    h6->SetMarkerColor(kBlack);
    h6->SetLineColor(kBlack);
    h6->DrawClone("ESAME");

    // Set the x-y axis labels and calculate the binning
    auto binning = (lim_right - lim_left)/bins;
    std::stringstream ss;
    ss.precision(2);
    ss << binning;
    std::string y_label = "Candidates / " + ss.str() + " " + unit;
    std::string x_label = xlabel_ + " " + unit;
    hs->GetXaxis()->SetTitle(x_label.c_str());
    hs->GetYaxis()->SetTitle(y_label.c_str());
    hs->SetTitle(title);

    // Draw a TLegend for Real Data and MC samples
    TLegend *legend = new TLegend(0.65, 0.65, 0.98, 0.88);
    legend->AddEntry(h1.GetPtr(),"gg#rightarrow H #rightarrow\#tau_{#mu}\#tau_{h}","f");
    legend->AddEntry(h2.GetPtr(),"q#rightarrow H #rightarrow\#tau_{#mu}\#tau_{h}","f");
    legend->AddEntry(h3.GetPtr(),"Z#rightarrow ll","f");
    legend->AddEntry(h4.GetPtr(),"t#bar{t}","f");
    legend->AddEntry(h5.GetPtr(),"W+jets","f");
    legend->AddEntry(h6.GetPtr(),"Data");
    legend->SetBorderSize(0);
    legend->SetFillStyle(0);
    legend->Draw();

    float xposition_lt = hs->GetXaxis()->GetXmax() *0.03 + hs->GetXaxis()->GetXmin();
    TLatex *lt = new TLatex(xposition_lt, hs->GetMaximum()*0.95,
        "#scale[0.8]{#splitline{CMS #bf{#it{Open Data}}}{#bf{11.5 fb^{-1} (8 TeV)}}}");
    lt->Draw("Same");

    // Save the plot in the correct directory
    std::string name = "./Distribution_plots/" + variable + ".pdf";
    c->SaveAs(name.c_str());
    c->Update();
}
```

# Plotting code

- Assegno i colori ad ogni istogramma e definisco un TStack per contenerli tutti

```
template <typename T>
void plotting(T samples, std::string variable, const char *title, std::string xlabel_,
std::string unit, int bins, float lim_left, float lim_right) {
    std::cout << "\tPlotting: " << variable << std::endl;

    TCanvas *c = new TCanvas("c","c", 1000, 700);
    gStyle->SetOptStat(0);
    gStyle->SetTitleFontSize(0.05);
    // Define six 1D histograms for each channel
    auto h1 = samples.at("GluGluToHToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h2 = samples.at("VBF_HToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h3 = samples.at("DYJetsToLL_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h4 = samples.at("TTbar_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h5 = samples.at("WJetsToLNu_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h6 = samples.at("Run2012*TauPlusX_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable);

    std::vector<decltype(h1)> hist = {h1, h2, h3, h4, h5};
    std::vector<decltype(kRed+1)> color =
        {kRed+1, kMagenta-9, kGreen-3, kAzure+7, kOrange-3};
    // Assign the colours to each histogram
    for (int i=0; i<hist.size(); i++) {
        hist[i]->SetLineColor(color[i]);
        hist[i]->SetFillColor(color[i]);
    }
    // Create a stacked histogram with MC simulation samples
    THStack *hs = new THStack("hs", "");
    for (auto &h : hist) {
        hs->Add(h.GetPtr());
    }
    hs->Draw("HIST0");
    // Draw the Data with black points and error bars
    h6->SetMarkerStyle(kFullCircle);
    h6->SetMarkerColor(kBlack);
    h6->SetLineColor(kBlack);
    h6->DrawClone("ESAME");

    // Set the x-y axis labels and calculate the binning
    auto binning = (lim_right - lim_left)/bins;
    std::stringstream ss;
    ss.precision(2);
    ss << binning;
    std::string y_label = "Candidates / " + ss.str() + " " + unit;
    std::string x_label = xlabel_ + " " + unit;
    hs->GetXaxis()->SetTitle(x_label.c_str());
    hs->GetYaxis()->SetTitle(y_label.c_str());
    hs->SetTitle(title);

    // Draw a Legend for Real Data and MC samples
    TLegend *legend = new TLegend(0.65, 0.65, 0.98, 0.88);
    legend->AddEntry(h1.GetPtr(),"gg#rightarrow H #rightarrow\#tau_{#mu}\#tau_{h}","f");
    legend->AddEntry(h2.GetPtr(),"q\bar{q}#rightarrow H #rightarrow\#tau_{#mu}\#tau_{h}","f");
    legend->AddEntry(h3.GetPtr(),"Z#rightarrow ll","f");
    legend->AddEntry(h4.GetPtr(),"t#bar{t}","f");
    legend->AddEntry(h5.GetPtr(),"W+jets","f");
    legend->AddEntry(h6.GetPtr(),"Data");
    legend->SetBorderSize(0);
    legend->SetFillStyle(0);
    legend->Draw();

    float xposition_lt = hs->GetXaxis()->GetXmax() *0.03 + hs->GetXaxis()->GetXmin();
    TLatex *lt = new TLatex(xposition_lt, hs->GetMaximum()*0.95,
        "#scale[0.8]{#splitline{CMS #bf{#it{Open Data}}}{#bf{11.5 fb^{-1} (8 TeV)}}}");
    lt->Draw("Same");

    // Save the plot in the correct directory
    std::string name = "./Distribution_plots/" + variable + ".pdf";
    c->SaveAs(name.c_str());
    c->Update();
}
```

# Plotting code

- Imposto le label per gli assi e creo la TLegend

```

template <typename T>
void plotting(T samples, std::string variable, const char *title, std::string xlabel_,
std::string unit, int bins, float lim_left, float lim_right) {
    std::cout << "\tPlotting: " << variable << std::endl;

    TCanvas *c = new TCanvas("c","c", 1000, 700);
    gStyle->SetOptStat(0);
    gStyle->SetTitleFontSize(0.05);
    // Define six 1D histograms for each channel
    auto h1 = samples.at("GluGluToHToTauTau_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable, "weight");
    auto h2 = samples.at("VBF_HToTauTau_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable, "weight");
    auto h3 = samples.at("DYJetsToLL_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable, "weight");
    auto h4 = samples.at("TTbar_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable, "weight");
    auto h5 = samples.at("WJetsToLNu_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable, "weight");
    auto h6 = samples.at("Run2012*TauPlusX_selected")
        .Histo1D({"", "", bins, lim_left, lim_right}, variable);

    std::vector<decltype(h1)> hist = {h1, h2, h3, h4, h5};
    std::vector<decltype(kRed+1)> color =
        {kRed+1, kMagenta-9, kGreen-3, kAzure+7, kOrange-3};
    // Assign the colours to each histogram
    for (int i=0; i<hist.size(); i++) {
        hist[i]->SetLineColor(color[i]);
        hist[i]->SetFillColor(color[i]);
    }
    // Create a stacked histogram with MC simulation samples
    THStack *hs = new THStack("hs", "");
    for (auto &h : hist) {
        hs->Add(h.GetPtr());
    }
    hs->Draw("HIST0");
    // Draw the Data with black points and error bars
    h6->SetMarkerStyle(kFullCircle);
    h6->SetMarkerColor(kBlack);
    h6->SetLineColor(kBlack);
    h6->DrawClone("ESAME");

    // Set the x-y axis labels and calculate the binning
    auto binning = (lim_right - lim_left)/bins;
    std::stringstream ss;
    ss.precision(2);
    ss << binning;
    std::string y_label = "Candidates / " + ss.str() + " " + unit;
    std::string x_label = xlabel_ + " " + unit;
    hs->GetXaxis()->SetTitle(x_label.c_str());
    hs->GetYaxis()->SetTitle(y_label.c_str());
    hs->SetTitle(title);

    // Draw a TLegend for Real Data and MC samples
    TLegend *legend = new TLegend(0.65, 0.65, 0.98, 0.88);
    legend->AddEntry(h1.GetPtr(),"gg#rightarrow H #rightarrow\tau_{#mu}\tau_{h}","f");
    legend->AddEntry(h2.GetPtr(),"q q#rightarrow H #rightarrow\tau_{#mu}\tau_{h}","f");
    legend->AddEntry(h3.GetPtr(),"Z#rightarrow ll","f");
    legend->AddEntry(h4.GetPtr(),"t#bar{t}","f");
    legend->AddEntry(h5.GetPtr(),"W+jets","f");
    legend->AddEntry(h6.GetPtr(),"Data");
    legend->SetBorderSize(0);
    legend->SetFillStyle(0);
    legend->Draw();

    float xposition_lt = hs->GetXaxis()->GetXmax()*0.03 + hs->GetXaxis()->GetXmin();
    TLatex *lt = new TLatex(xposition_lt, hs->GetMaximum()*0.95,
        "#scale[0.8]{#splitline{CMS #bf{#it{Open Data}}}{#bf{11.5 fb^{-1} (8 TeV)}}}");
    lt->Draw("Same");

    // Save the plot in the correct directory
    std::string name = "./Distribution_plots/" + variable + ".pdf";
    c->SaveAs(name.c_str());
    c->Update();
}

```

# Plotting code

- Salvo tutti i plot nel folder “Distribution Plots”

```

template <typename T>
void plotting(T samples, std::string variable, const char *title, std::string xlabel_,
std::string unit, int bins, float lim_left, float lim_right) {
    std::cout << "\tPlotting: " << variable << std::endl;

    TCanvas *c = new TCanvas("c","c", 1000, 700);
    gStyle->SetOptStat(0);
    gStyle->SetTitleFontSize(0.05);
    // Define six 1D histograms for each channel
    auto h1 = samples.at("GluGluToHToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h2 = samples.at("VBF_HToTauTau_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h3 = samples.at("DYJetsToLL_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h4 = samples.at("TTbar_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h5 = samples.at("WJetsToLNu_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable, "weight");
    auto h6 = samples.at("Run2012*TauPlusX_selected")
        .Histo1D({ "", "", bins, lim_left, lim_right}, variable);

    std::vector<decltype(h1)> hist = {h1, h2, h3, h4, h5};
    std::vector<decltype(kRed+1)> color =
        {kRed+1, kMagenta-9, kGreen-3, kAzure+7, kOrange-3};
    // Assign the colours to each histogram
    for (int i=0; i<hist.size(); i++) {
        hist[i]->SetLineColor(color[i]);
        hist[i]->SetFillColor(color[i]);
    }
    // Create a stacked histogram with MC simulation samples
    THStack *hs = new THStack("hs", "");
    for (auto &h : hist) {
        hs->Add(h.GetPtr());
    }
    hs->Draw("HIST0");
    // Draw the Data with black points and error bars
    h6->SetMarkerStyle(kFullCircle);
    h6->SetMarkerColor(kBlack);
    h6->SetLineColor(kBlack);
    h6->DrawClone("ESAME");

    // Set the x-y axis labels and calculate the binning
    auto binning = (lim_right - lim_left)/bins;
    std::stringstream ss;
    ss.precision(2);
    ss << binning;
    std::string y_label = "Candidates / " + ss.str() + " " + unit;
    std::string x_label = xlabel_ + " " + unit;
    hs->GetXaxis()->SetTitle(x_label.c_str());
    hs->GetYaxis()->SetTitle(y_label.c_str());
    hs->SetTitle(title);

    // Draw a TLegend for Real Data and MC samples
    TLegend *legend = new TLegend(0.65, 0.65, 0.98, 0.88);
    legend->AddEntry(h1.GetPtr(),"gg#rightarrow H #rightarrow\tau_{#mu}\tau_{h}","f");
    legend->AddEntry(h2.GetPtr(),"q q#rightarrow H #rightarrow\tau_{#mu}\tau_{h}","f");
    legend->AddEntry(h3.GetPtr(),"Z#rightarrow ll","f");
    legend->AddEntry(h4.GetPtr(),"t#bar{t}","f");
    legend->AddEntry(h5.GetPtr(),"W+jets","f");
    legend->AddEntry(h6.GetPtr(),"Data");
    legend->SetBorderSize(0);
    legend->SetFillStyle(0);
    legend->Draw();

    float xposition_lt = hs->GetXaxis()->GetXmax() *0.03 + hs->GetXaxis()->GetXmin();
    TLatex *lt = new TLatex(xposition_lt, hs->GetMaximum()*0.95,
        "#scale[0.8]{#splitline{CMS #bf{#it{Open Data}}}{#bf{11.5 fb^{-1} (8 TeV)}}}");
    lt->Draw("Same");

    // Save the plot in the correct directory
    std::string name = "./Distribution_plots/" + variable + ".pdf";
    c->SaveAs(name.c_str());
    c->Update();
}

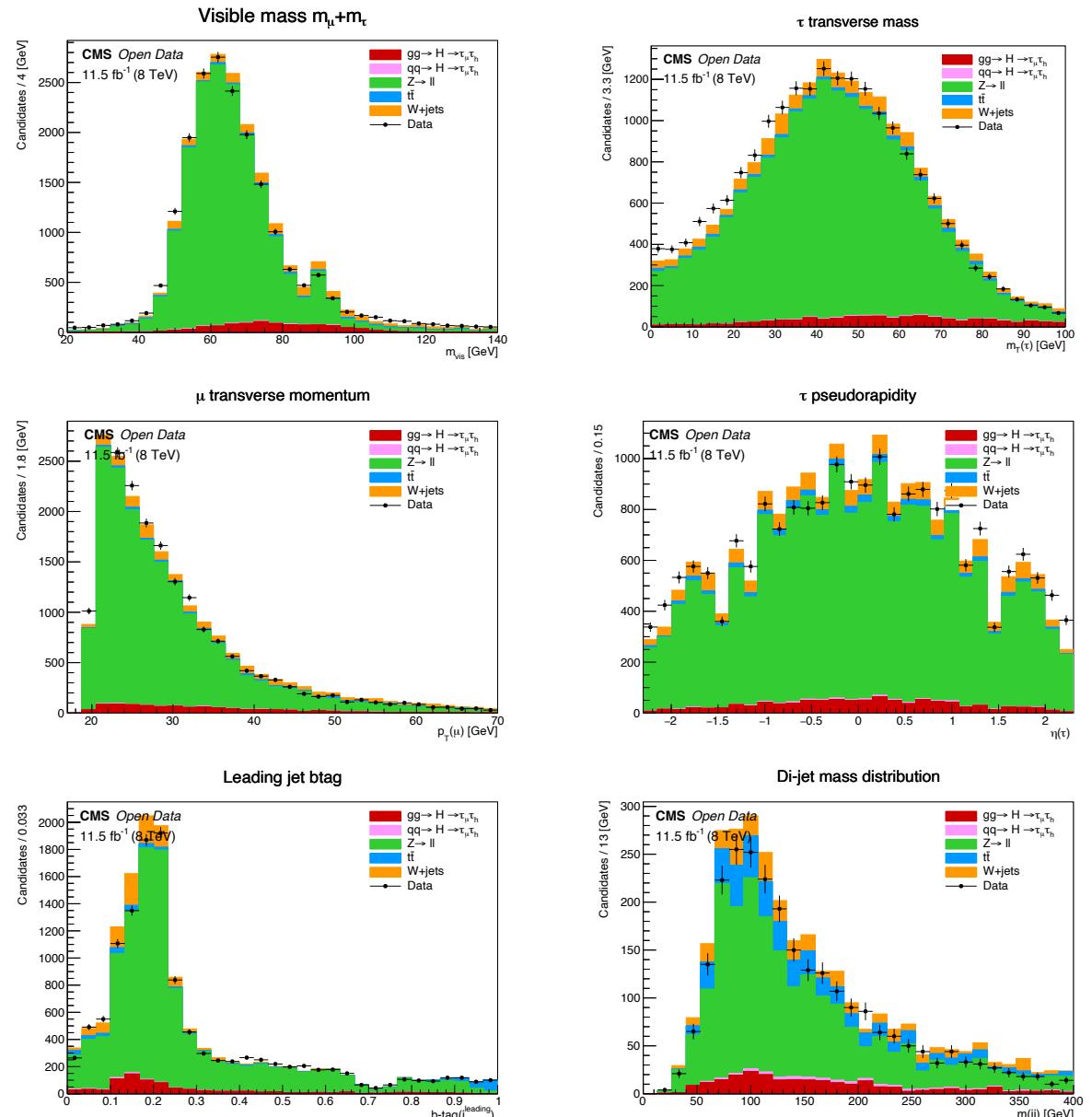
```

# Plotting code

```

domenico@MacBook-Pro-di-domenico ROOT_workspace % ./plotting
Processing: GluGluToHToTauTau_selected
Number of events: 2034
Processing: VBF_HToTauTau_selected
Number of events: 2514
Processing: DYJetsToLL_selected
Number of events: 12927
Processing: TTbar_selected
Number of events: 1344
Processing: W*JetsToLNu_selected
Number of events: 1410
Processing: Run2012*_TauPlusX_selected
Number of events: 20187
***** Plotting distribution *****
Plotting: m_vis
Plotting: pt_vis
Plotting: tau_pt
Plotting: muon_pt
Plotting: muon_eta
Plotting: tau_eta
Plotting: muon_phi
Plotting: tau_phi
Plotting: muon_iso
Plotting: tau_iso
Plotting: muon_charge
Plotting: tau_charge
Plotting: MET_pt
Plotting: MET_phi
Plotting: muon_m
Plotting: tau_m
Plotting: mt_mu
Plotting: mt_tau
Plotting: mt_mu
Plotting: jpt_1
Plotting: jpt_2
Plotting: jeta_1
Plotting: jeta_2
Plotting: jphi_1
Plotting: jphi_2
Plotting: jm_1
Plotting: jm_2
Plotting: jbtags_1
Plotting: jbtags_2
Plotting: PV_npvs
Plotting: nGoodJets
Plotting: jj_m
Plotting: jj_pt
Plotting: jj_delta

```



Distribution\_plots

# **Algoritmi ML per la classificazione segnale/fondo (Python)**

Seconda Parte



# Machine Learning: classificazione segnale/fondo (1)

## Main function

```
1 if __name__ == "__main__":
2     ML_dict = { ... }
3
4     # Create a new directory contains ML plots
5     NAME_DIR = '/ML_plots'
6     PATH = os.getcwd()
7     if not os.path.exists(PATH + NAME_DIR):
8         os.mkdir(PATH + NAME_DIR)
9         print(f">>> ..{NAME_DIR} directory has been created at {PATH}"")
10    else:
11        print(f">>> ..{NAME_DIR} folder is already present at {PATH}"")
12
13    # Ask to user what signal channel would like to use
14    print("Which file would you process? [ggH o VBF]")
15    CHOICE = input()
16    print(f"*****{CHOICE} is our signal channel *****")
17    if (CHOICE != "ggH") and (CHOICE != "VBF"):
18        print("This CHOICE isn't present")
19        sys.exit()
20
21    events = read_root(ML_dict[CHOICE]['complete_name'])
22    NUM_VARS = len(ML_dict[CHOICE]['ML_VARS'])
23    # Remove undefined variables
24    for i in range(NUM_VARS):
25        events = events[(events[VARS[i]] > -999)]
26
27    print(f"Number of signal (1) and background events (0):\n{events.event.value_counts()}"")
```

```
ML_dict = {
    'ggH': {
        'complete_name': 'GluGluToHToTauTau',
        'ML_VARS': ["nGoodJets", "PV_npvs", "MET_pt", "MET_phi", "pt_vis",
                    "muon_pt", "muon_eta", "muon_phi", "muon_iso",
                    "tau_pt", "tau_eta", "tau_phi"],
        'number_input': 10,
        'N_EPOCHS': 150,
    },
    'VBF': {
        'complete_name': 'VBF_HToTauTau',
        'ML_VARS': ["nGoodJets", "PV_npvs", "MET_phi", "m_vis", "pt_vis",
                    "muon_pt", "muon_eta", "muon_phi", "muon_iso",
                    "tau_pt", "tau_eta", "tau_phi", "tau_iso",
                    "jbtag_1", "jj_delta"],
        'number_input': 15,
        'N_EPOCHS': 100
    }
}
```

# Machine Learning: classificazione segnale/fondo (2)

## Read root function

```
1 def read_root(file_name):
2     """ doc """
3     PATH = "ROOT_workspace/"
4     df = {}
5     for file in MLfiles:
6         Events = uproot.open(PATH + file + "_selected.root" + ":Events")
7         df[file] = pd.DataFrame(Events.arrays(VARS, library="np"), columns=VARS)
8         print(f"\tProcessing: {file}, Number of events: {len(df[file])}")
9         if file == file_name:
10             df[file]["event"] = 1.
11         else:
12             df[file]["event"] = 0.
13     return pd.concat(df.values(), ignore_index=True)
```

```
● ● ●
*****
ggH is our signal channel *****
Processing: GluGluToHToTauTau, Number of events: 4492
Processing: VBF_HToTauTau, Number of events: 5326
Processing: DYJetsToLL, Number of events: 35099
Processing: TTbar, Number of events: 13630
Processing: W1JetsToLNu, Number of events: 7042
Processing: W2JetsToLNu, Number of events: 11427
Processing: W3JetsToLNu, Number of events: 6327
Number of signal (1) and background events (0):
0.0    78839
1.0    4492
Name: event, dtype: int64
```

# Machine Learning: classificazione segnale/fondo (3)

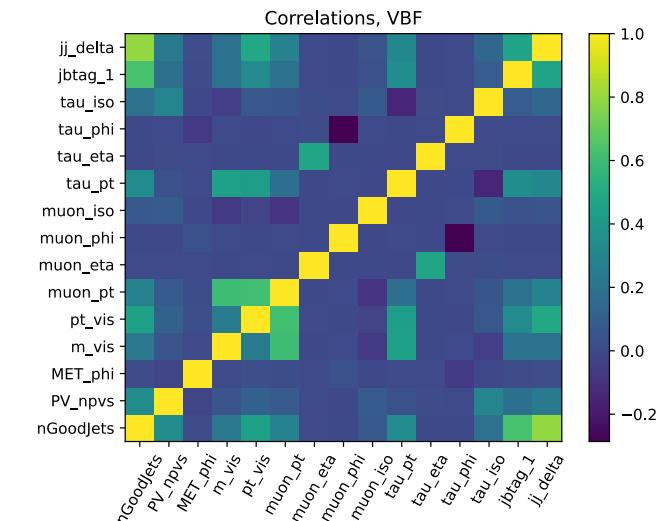
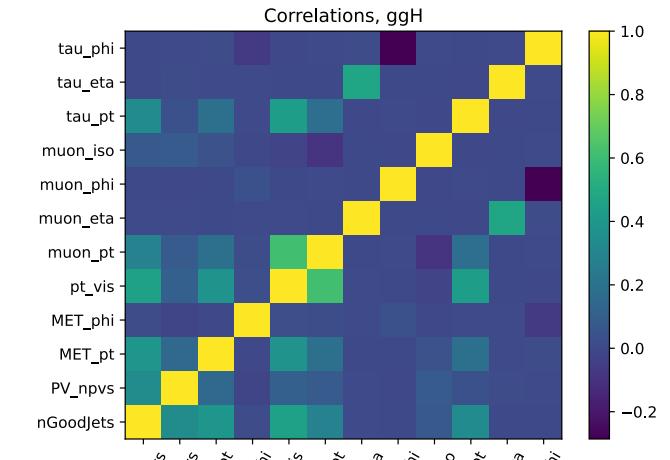
## Main function

```
1 # Renormalization of events
2 sig = events["event"] == 1.
3 bkg = events["event"] == 0.
4 weight_sig_sum = events["weight"][sig].sum(axis=0)
5 weight_bkg_sum = events["weight"][bkg].sum(axis=0)
6 events["weight"][sig] = events["weight"][sig] / weight_sig_sum
7 events["weight"][bkg] = events["weight"][bkg] / weight_bkg_sum
8
9 # Plot correlation matrix of ML variables
10 correlations(CHOICE, events[ML_dict[CHOICE][ "ML_VARS"]])
11
12 # Generate features' matrix
13 features = events.filter(ML_dict[CHOICE][ "ML_VARS"])
14 X = np.asarray(features.values).astype(np.float32)
15 # Generate target vector
16 target = events.filter(['event'])
17 y = np.asarray(target.values).astype(np.float32)
18 # Generate event weights vector
19 weight = events.filter(['weight'])
20 W = np.asarray(weight.values).astype(np.float32)
21
22 # Split into training and testing dataset
23 X_train, X_test, y_train, y_test, W_train, W_test = train_test_split(X, y, W,
24 test_size=0.3, shuffle=True)
25 print(f'Number of events for training phase: {len(X_train)}')
26 print(f'Number of events for test phase: {len(X_test)}')
27
28 # Perform features scaling with StandardScaler
29 SC = StandardScaler()
30 X_train = SC.fit_transform(X_train)
31 X_test = SC.transform(X_test)
```

# Machine Learning: classificazione segnale/fondo (4)

## Correlation function

```
1 def correlations(channel_name, data):
2     """ doc """
3     plt.figure()
4     corrmat = data.corr()
5     heatmap1 = plt.pcolor(corrmat) # get heatmap
6     plt.colorbar(heatmap1) # plot colorbar
7
8     plt.title(f'Correlations, {channel_name}') # set title
9     x_variables = corrmat.columns.values # get variables from data columns
10
11    plt.xticks(np.arange(len(x_variables)) + 0.5, x_variables, rotation=60) # x-tick for each label
12    plt.yticks(np.arange(len(x_variables)) + 0.5, x_variables) # y-tick for each label
13    plt.savefig(f'{NAME_DIR}/correlation_{channel_name}.pdf', format='pdf', bbox_inches='tight')
14    plt.clf()
```



# ANN Machine Learning: classificazione segnale/fondo (1)

## Main function, ANN

```
1  print("\n***** ARTIFICIAL NEURAL NETWORK *****")
2  print(f"\tML variables ({NUM_VARS} features) are:")
3  for var in ML_dict[CHOICE]['ML_VARS']:
4      print(f"\t{var}")
5
6  # Create Artificial Neural Network (with 3 hidden layers)
7  input_layer = Input(shape=(NUM_VARS,), name='input')
8  hidden = Dense(NUM_VARS * ML_dict[CHOICE]['number_input'], name='hidden1',
9                 activation='selu')(input_layer)
10 hidden = Dropout(rate=0.1)(hidden)
11 hidden = Dense(NUM_VARS * 2, name='hidden2', activation='selu')(hidden)
12 hidden = Dropout(rate=0.5)(hidden)
13 hidden = Dense(NUM_VARS, name='hidden3', activation='selu')(hidden)
14 hidden = Dropout(rate=0.1)(hidden)
15 output_layer = Dense(1, name='output', activation='sigmoid')(hidden)
16 # Initialising ANN model
17 ANN = Model(inputs=input_layer, outputs=output_layer, name='ANN')
18 ANN.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'],
19             weighted_metrics=['accuracy']) # Compiling ANN
20 ANN.summary() # Print on screen the model summary
21
22 # Fit the ANN model
23 HISTORY = ANN.fit(X_train, y_train, sample_weight=W_train, batch_size=20,
24                     epochs=ML_dict[CHOICE]["N_EPOCHS"], verbose=1, validation_split=0.3)
25
26 plotting_loss(CHOICE, HISTORY.history['loss'], HISTORY.history['val_loss'])
27 plotting_accuracy(CHOICE, HISTORY.history['accuracy'], HISTORY.history['val_accuracy'])
```

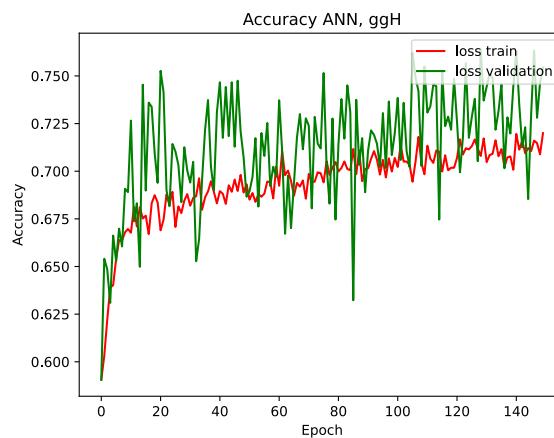
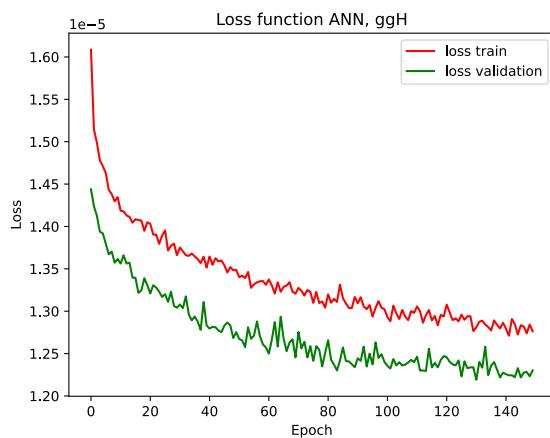
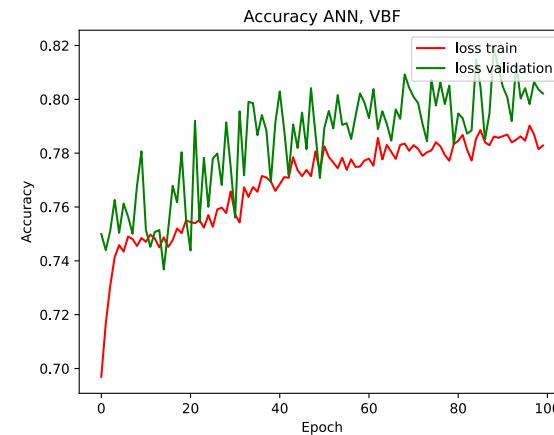
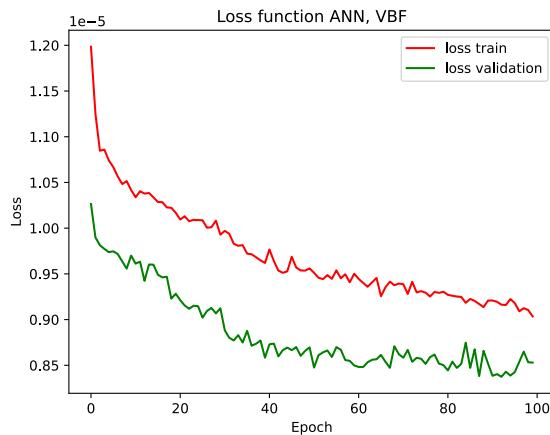
ggH		
Layer (type)	Output Shape	Param #
input (InputLayer)	(None, 12)	0
hidden1 (Dense)	(None, 120)	1560
dropout (Dropout)	(None, 120)	0
hidden2 (Dense)	(None, 24)	2904
dropout_1 (Dropout)	(None, 24)	0
hidden3 (Dense)	(None, 12)	300
dropout_2 (Dropout)	(None, 12)	0
output (Dense)	(None, 1)	13
Total params: 4,777		
Trainable params: 4,777		
Non-trainable params: 0		

# ANN Machine Learning: classificazione segnale/fondo (2)

## Loss and accuracy plotting functions

```
1 def plotting_loss(channel_name, loss, val_loss):
2     """ doc """
3     plt.figure()
4     plt.plot(loss, label='loss train', color='red')
5     plt.plot(val_loss, label='loss validation', color='green')
6     plt.title(f'Loss function ANN, {channel_name}')
7     plt.xlabel('Epoch')
8     plt.ylabel('Loss')
9     plt.legend(loc='upper right')
10    plt.savefig(f'{NAME_DIR}/Loss_ANN_{channel_name}.pdf',
11                format='pdf')
12    plt.clf()
13
14 def plotting_accuracy(channel_name, accuracy, val_accuracy):
15     """ doc """
16     plt.figure()
17     plt.plot(accuracy, label='loss train', color='red')
18     plt.plot(val_accuracy, label='loss validation', color='green')
19     plt.title(f'Accuracy ANN, {channel_name}')
20     plt.xlabel('Epoch')
21     plt.ylabel('Accuracy')
22     plt.legend(loc='upper right')
23     plt.savefig(f'{NAME_DIR}/Accuracy_ANN_{channel_name}.pdf',
24                 format='pdf')
25     plt.clf()
```

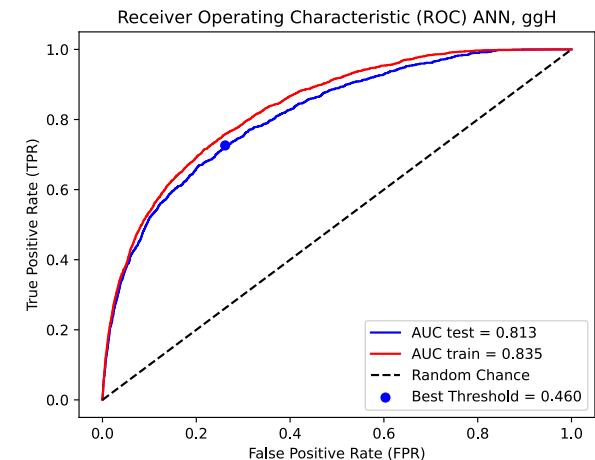
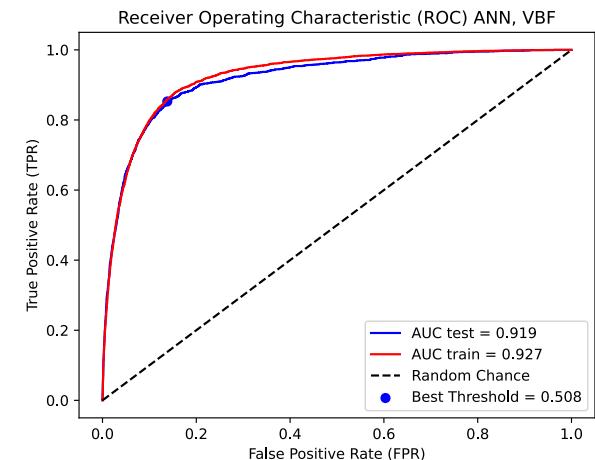
# ANN Machine Learning: classificazione segnale/fondo (3)



# ANN Machine Learning: classificazione segnale/fondo (4)

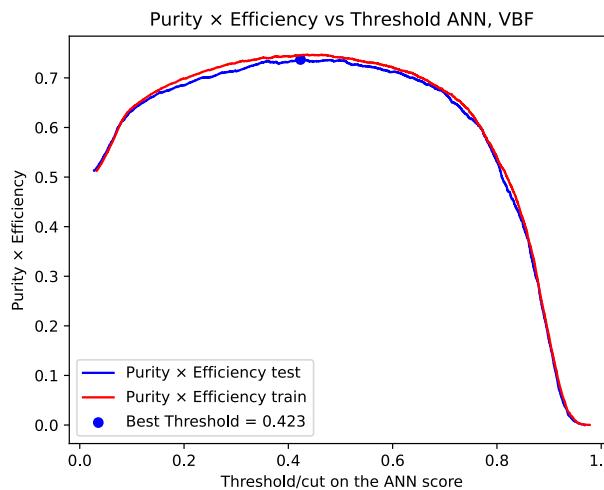
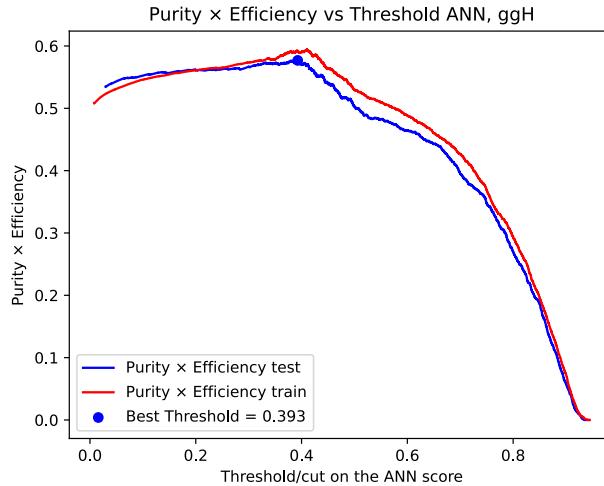
## Main function, ROC curves

```
1 # Get ANN model event predictions
2 y_prediction_test = ANN.predict(X_test)
3 y_prediction_train = ANN.predict(X_train)
4
5 # Get False Positive Rate (FPR), True Positive Rate (TPR) and Thresholds/Cut
6 # on the ANN's score
7 fpr_test, tpr_test, thresholds_test = roc_curve(y_true=y_test,
8                                                 y_score=y_prediction_test,
9                                                 sample_weight=W_test)
10 fpr_train, tpr_train, thresholds_train = roc_curve(y_true=y_train,
11                                                 y_score=y_prediction_train,
12                                                 sample_weight=W_train)
13
14 # Optimal Threshold for ROC Curve with g-mean
15 index_g = g_mean_threshold(tpr_test, fpr_test, thresholds_test)
16
17
18 # Plot the ANN ROC curve on the test and training datasets
19 # with plotting_ROC function implemented in plot_function.py script
20 roc_auc_test = auc(fpr_test, tpr_test)
21 roc_auc_train = auc(fpr_train, tpr_train)
22 plotting_ROC(f'ANN_{CHOICE}', fpr_test, tpr_test, fpr_train, tpr_train,
23               thresholds_test, index_g, roc_auc_test, roc_auc_train)
```



# ANN Machine Learning: classificazione segnale/fondo (5)

## Main function, P-E curves



```
1 # Get precision (purity) and recall (signal efficiency) with precision_recall_curve method.
2 purity_test, recall_test, t_test = precision_recall_curve(y_true=y_test,
3                                                               probas_pred=y_prediction_test,
4                                                               sample_weight=w_test)
5 purity_train, recall_train, t_train = precision_recall_curve(y_true=y_train,
6                                                               probas_pred=y_prediction_train,
7                                                               sample_weight=w_train)
8
9 # Optimal threshold for Precision-Recall Curve (PR) with simplify f-score
10 index_f = fscore_threshold(purity_test[:-1], recall_test[:-1], t_test)
11
12 plotting_purity_vs_efficiency(f'ANN_{CHOICE}', t_test, t_train, recall_test,
13                                recall_train, purity_test, purity_train, index_f)
```

# ANN Machine Learning: classificazione segnale/fondo (6)

## Plotting function, ROC and P-E

```
1 def plotting_ROC(channel_name, fpr_test, tpr_test, fpr_train, tpr_train, thresholds,
2                   ix, roc_auc_test, roc_auc_train):
3     """ doc """
4     names = channel_name.split("_")
5     plt.figure()
6     plt.plot(fpr_test, tpr_test, color='blue', label=f'AUC test = {roc_auc_test:.3f}')
7     plt.plot(fpr_train, tpr_train, color='red', label=f'AUC train = {roc_auc_train:.3f}')
8     plt.plot([0, 1], [0, 1], linestyle='--', color='k', label='Random Chance')
9     plt.scatter(fpr_test[ix], tpr_test[ix], marker='o', color='blue',
10                 label=f'Best Threshold = {thresholds[ix]:.3f}')
11    plt.xlabel('False Positive Rate (FPR)')
12    plt.ylabel('True Positive Rate (TPR)')
13    plt.title(f'Receiver Operating Characteristic (ROC) {names[0]}, {names[1]}')
14    plt.legend(loc='lower right')
15    plt.savefig(f'{NAME_DIR}/ROC_{channel_name}.pdf', format='pdf')
16    plt.clf()
17
18 def plotting_purity_vs_efficiency(channel_name, t_test, t_train, recall_test, recall_train,
19                                    purity_test, purity_train, ix):
20     """ doc """
21     names = channel_name.split("_")
22     plt.figure()
23     plt.plot(t_test, recall_test[:-1] * purity_test[:-1], color='blue',
24               label=r'Purity $\times$ Efficiency test')
25     plt.plot(t_train, recall_train[:-1] * purity_train[:-1], color='red',
26               label=r'Purity $\times$ Efficiency train')
27     plt.scatter(t_test[ix], recall_test[ix] * purity_test[ix], marker='o', color='blue',
28                 label=f'Best Threshold = {t_test[ix]:.3f}')
29     plt.xlabel(f'Threshold/cut on the {names[0]} score')
30     plt.ylabel(r'Purity $\times$ Efficiency')
31     plt.title(f'Purity $\times$ Efficiency vs Threshold {names[0]}, {names[1]}')
32     plt.legend(loc='lower left')
33     plt.savefig(f'{NAME_DIR}/metrics_PR_{channel_name}.pdf', format='pdf')
34     plt.clf()
```

# ANN Machine Learning: classificazione segnale/fondo (7)

## Main function, g-mean and simply f-score

```
1 def g_mean_threshold(tpr, fpr, thresholds):
2     """ doc """
3     g_means = np.sqrt(tpr * (1 - fpr))
4     index = np.argmax(g_means) # locate the index of the largest g-mean
5     print(f'Best Threshold={thresholds[index]:.3f}, \
6 G-Mean={g_means[index]:.3f}')
7     return index
8
9 def fscore_threshold(purity, recall, thresholds):
10    """ doc """
11    fscore = purity[:-1] * recall[:-1]
12    index = np.argmax(fscore) # locate the index of the largest f-score
13    print(f'Best Threshold={thresholds[index]:.3f}, \
14 F-Score simplify={fscore[index]:.3f}')
15    return index
```

# ANN Machine Learning: classificazione segnale/fondo (8)

## Main function, confusion matrix and res.

```
1 # Threshold selection according to signal channel
2 cut_dnn = 0.
3 if CHOICE == "ggH":
4     cut_dnn = thresholds_test[index_g]
5 else:
6     cut_dnn = t_test[index_f]
7
8 # Plotting output score distribution
9 plotting_output_score(f'ANN_{CHOICE}', y_train, y_prediction_train, y_test,
10                      y_prediction_test, cut_dnn)
11
12 # Transform predictions into an array of 0, 1
13 filter_test_sig = y_prediction_test >= cut_dnn # classify as signal
14 filter_test_bkg = y_prediction_test < cut_dnn # classify as background
15 y_prediction_test[filter_test_sig] = 1
16 y_prediction_test[filter_test_bkg] = 0
17
18 # Metrics values for the ANN algorithm having fixed an ANN score threshold
19 accuracy = accuracy_score(y_test[:, 0], y_prediction_test[:, 0], sample_weight=W_test[:, 0])
20 precision = precision_score(y_test[:, 0], y_prediction_test[:, 0], sample_weight=W_test[:, 0])
21 recall = recall_score(y_test[:, 0], y_prediction_test[:, 0], sample_weight=W_test[:, 0])
22 print(f'Threshold on the ANN output : {cut_dnn:.3f}')
23 print(f'ANN Test Accuracy: {accuracy:.3f}')
24 print(f'ANN Test Precision/Purity: {precision:.3f}')
25 print(f'ANN Test Sensitivity/Recall/TPR/Signal Efficiency: {recall:.3f}')
26 plotting_confusion_matrix(f'ANN_{CHOICE}', y_test[:, 0], y_prediction_test[:, 0], W_test[:, 0])
```

```
●●● ggH
...
Epoch 150/150
2042/2042 [=====] - 5s 3ms/step -
loss: 1.2763e-05 - accuracy: 0.7200 - weighted_accuracy: 0.7248 -
val_loss: 1.2301e-05 - val_accuracy: 0.7544 - val_weighted_accuracy: 0.7356

Best Threshold=0.460, G-Mean=0.732
Best Threshold=0.393, F-Score simplify=0.577
Threshold on the ANN output : 0.460
ANN Test Accuracy: 0.732
ANN Test Precision/Purity: 0.741
ANN Test Sensitivity/Recall/TPR/Signal Efficiency: 0.726
```

```
●●● VBF
...
Epoch 100/100
2042/2042 [=====] - 13s 6ms/step -
loss: 9.0350e-06 - accuracy: 0.7829 - weighted_accuracy: 0.8481 -
val_loss: 8.5298e-06 - val_accuracy: 0.8022 - val_weighted_accuracy: 0.8594

Best Threshold=0.508, G-Mean=0.857
Best Threshold=0.423, F-Score simplify=0.737
Threshold on the ANN output : 0.423
ANN Test Accuracy: 0.854
ANN Test Precision/Purity: 0.840
ANN Test Sensitivity/Recall/TPR/Signal Efficiency: 0.877
```

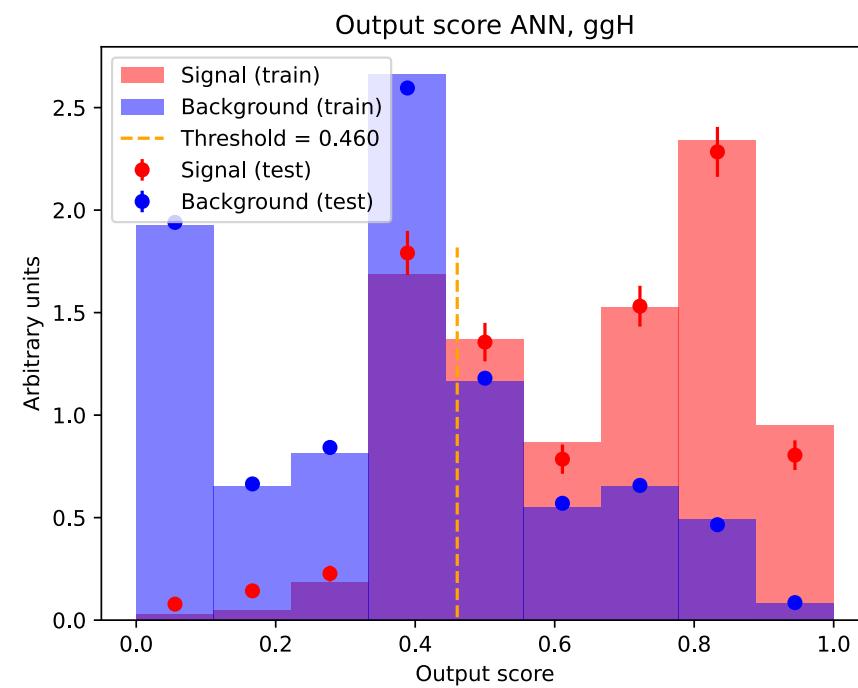
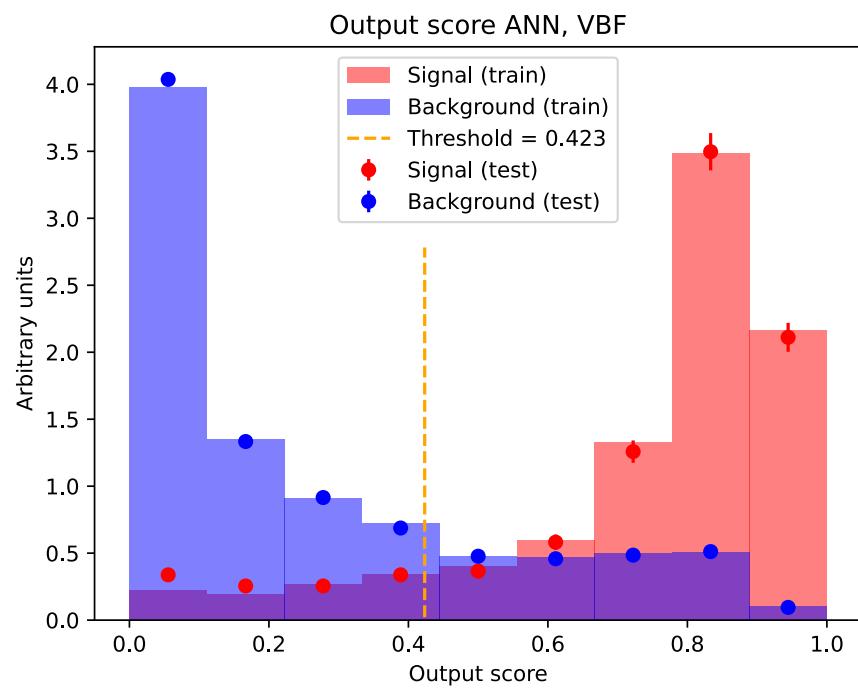
# ANN Machine Learning: classificazione segnale/fondo (9)

## Plotting function, output-score

```
1 def plotting_output_score(channel_name, y_train, y_prediction_train, y_test, y_prediction_test, cut_dnn):
2     """ doc """
3     names = channel_name.split("_")
4     plt.figure()
5     xx = np.linspace(0.0, 1.0, 10)
6     plt.hist(y_prediction_train[y_train == 1], bins=xx, density=1, label='Signal (train)',
7               alpha=0.5, color='red')
8     plt.hist(y_prediction_train[y_train == 0], bins=xx, density=1, label='Background (train)',
9               alpha=0.5, color='blue')
10
11    hist_test_sig, bin_edges = np.histogram(y_prediction_test[y_test == 1], bins=xx, density=True)
12    hist_test_bgk, _ = np.histogram(y_prediction_test[y_test == 0], bins=xx, density=True)
13
14    center = (bin_edges[:-1] + bin_edges[1:]) / 2 # bin centres
15
16    scale = len(y_prediction_test[y_test == 1]) / sum(hist_test_sig)
17    err_sig = np.sqrt(hist_test_sig * scale) / scale
18    plt.errorbar(center, hist_test_sig, yerr=err_sig, label='Signal (test)', c='red', fmt='o')
19
20    scale = len(y_prediction_test[y_test == 0]) / sum(hist_test_bgk)
21    err_bgk = np.sqrt(hist_test_bgk * scale) / scale
22    plt.errorbar(center, hist_test_bgk, yerr=err_bgk, label='Background (test)', c='blue', fmt='o')
23    plt.axvline(x=cut_dnn, ymin=0.005, ymax=0.65, color='orange', label=f'Threshold = {cut_dnn:.3f}',
24                 linestyle='--')
25    plt.title(f'Output score {names[0]}, {names[1]}')
26    plt.xlabel('Output score')
27    plt.ylabel('Arbitrary units')
28    plt.legend()
29    plt.savefig(f'{NAME_DIR}/{channel_name}_output_score.pdf', format='pdf')
30    plt.clf()
```

# ANN Machine Learning: classificazione segnale/fondo (10)

## Output score ANN

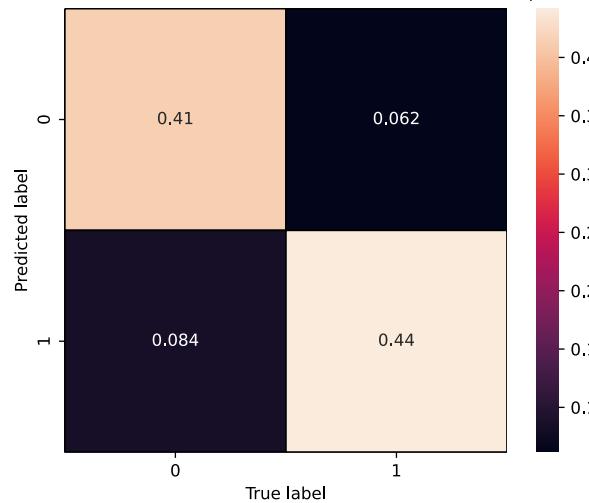


# ANN Machine Learning: classificazione segnale/fondo (11)

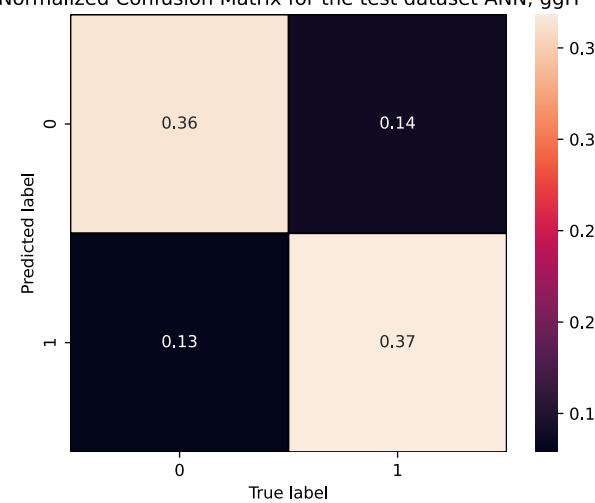
## Plotting function, confusion matrix

```
1 def plotting_confusion_matrix(channel_name, y_test, y_prediction_test, w_test):
2     """ doc """
3     names = channel_name.split("_")
4     plt.figure()
5     mat_test = confusion_matrix(y_test, y_prediction_test, sample_weight=w_test, normalize='all')
6     sns.heatmap(mat_test.T, square=True, annot=True, cbar=True, linewidths=1, linecolor='black')
7     plt.xlabel('True label')
8     plt.ylabel('Predicted label')
9     plt.title(f'Normalized Confusion Matrix for the test dataset {names[0]}, {names[1]}')
10    plt.savefig(f'{NAME_DIR}/confusion_matrix_{channel_name}.pdf', format='pdf')
11    plt.clf()
```

Normalized Confusion Matrix for the test dataset ANN, VBF



Normalized Confusion Matrix for the test dataset ANN, ggH



# RF Machine Learning: classificazione segnale/fondo (1)

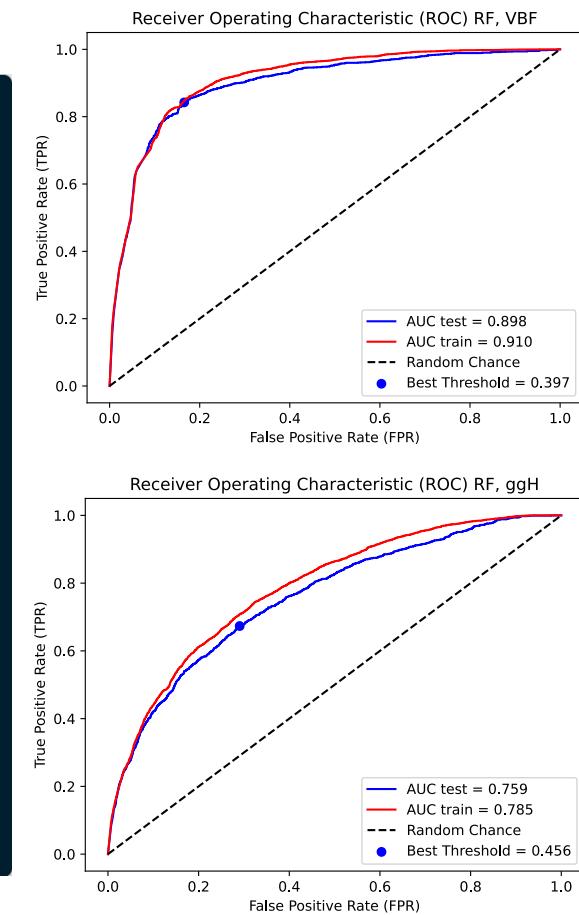
## Main function, model e Hyperparameter tuning

```
1 print("***** RANDOM FOREST *****")
2
3 # Uncomment the following lines to perform the tuning of hyper-parameters of the Random Forest Classifier
4 # classifier = RandomForestClassifier(random_state=7, verbose=1)
5 # grid_param = {
6 #     'criterion': ['gini', 'entropy'],
7 #     'n_estimators': [300, 500],
8 #     'bootstrap': [True, False],
9 #     'max_depth': [5, 7],
10 #     'max_features': [10, None],
11 # }
12 # tuner_RF = GridSearchCV(estimator=classifier, param_grid=grid_param, scoring='accuracy', cv=3, n_jobs=-1)
13 # tuner_RF.fit(X_train, np.ravel(y_train))
14 # print(f"Best parameters: {tuner_RF.best_params_}")
15 # print(f"Best metrics score {tuner_RF.best_score_}")
16
17 # Create a RF classifier
18 RFC = RandomForestClassifier(n_estimators=500, criterion='gini', verbose=1, max_depth=5,
19                             max_features=None, bootstrap=True, n_jobs=-1)
20
21 # Train the model on the training dataset
22 randomforest = RFC.fit(X_train, np.ravel(y_train), np.ravel(W_train))
23
24 # Get RF model event predictions on the test and training dataset.
25 # Return a vector of 0s and 1s (the two labels of the classification problem)
26 y_prediction_test_rf = RFC.predict(X_test)
27 y_prediction_train_rf = RFC.predict(X_train)
```

# RF Machine Learning: classificazione segnale/fondo (2)

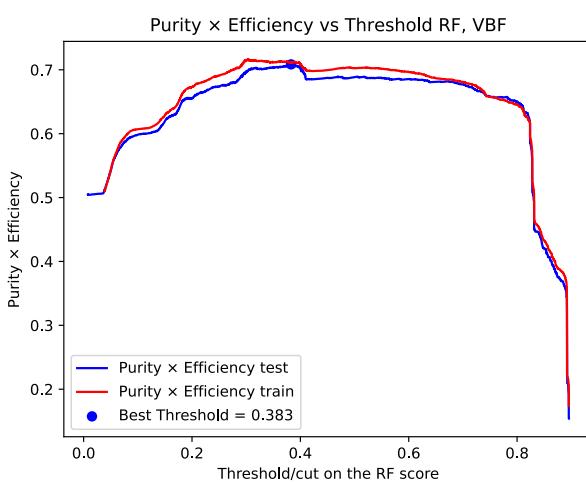
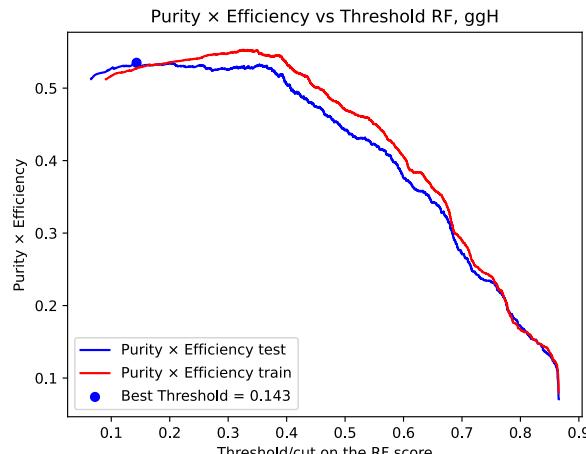
## Main function, ROC curves

```
1 # Method predict_proba return a 2D vector.
2 # The first column is the probability of the event being 'bkg'
3 # and the second is the probability of it being 'sig'
4 y_prediction_test_prob = randomforest.predict_proba(X_test)
5 y_prediction_train_prob = randomforest.predict_proba(X_train)
6
7 print(f"Accuracy of the RF model: {accuracy_score(y_test, y_prediction_test_rf):.3f} (test dataset)")
8 print(f"Accuracy of the RF model: {accuracy_score(y_train, y_prediction_train_rf):.3f} (train dataset)")
9
10 fpr_test_rf, tpr_test_rf, thresholds_test_rf = roc_curve(y_true=y_test,
11                                         y_score=y_prediction_test_prob[:, -1],
12                                         sample_weight=W_test)
13 fpr_train_rf, tpr_train_rf, thresholds_train_rf = roc_curve(y_true=y_train,
14                                         y_score=y_prediction_train_prob[:, -1],
15                                         sample_weight=W_train)
16
17 # Calculate area under the ROC curve
18 roc_auc_test_rf = auc(fpr_test_rf, tpr_test_rf)
19 roc_auc_train_rf = auc(fpr_train_rf, tpr_train_rf)
20
21 index_g_RF = g_mean_threshold(tpr_test_rf, fpr_test_rf, thresholds_test_rf)
22 plotting_ROC(f'RF_{CHOICE}', fpr_test_rf, tpr_test_rf, fpr_train_rf, tpr_train_rf,
23               thresholds_test_rf, index_g_RF, roc_auc_test_rf, roc_auc_train_rf)
24
```



# RF Machine Learning: classificazione segnale/fondo (3)

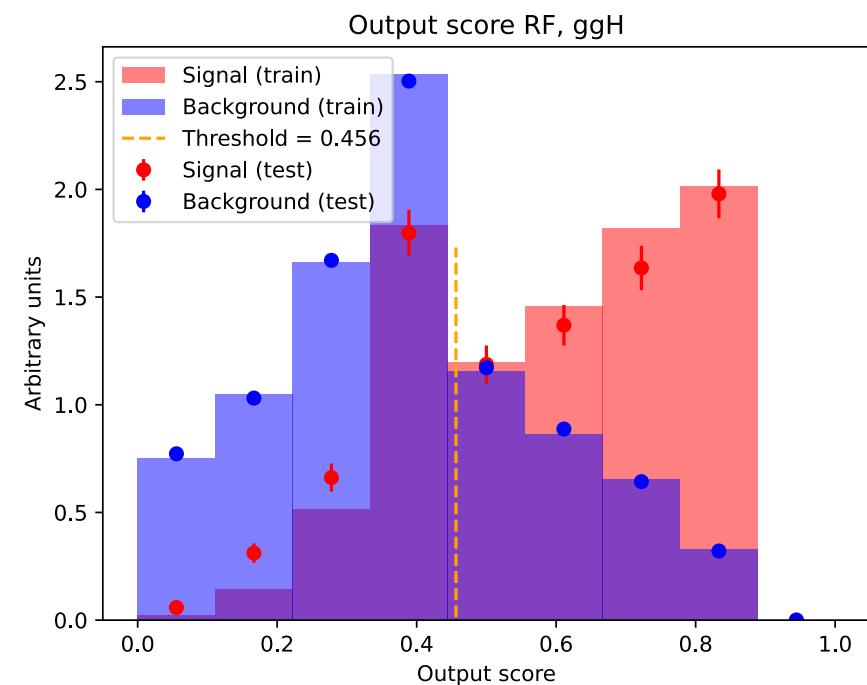
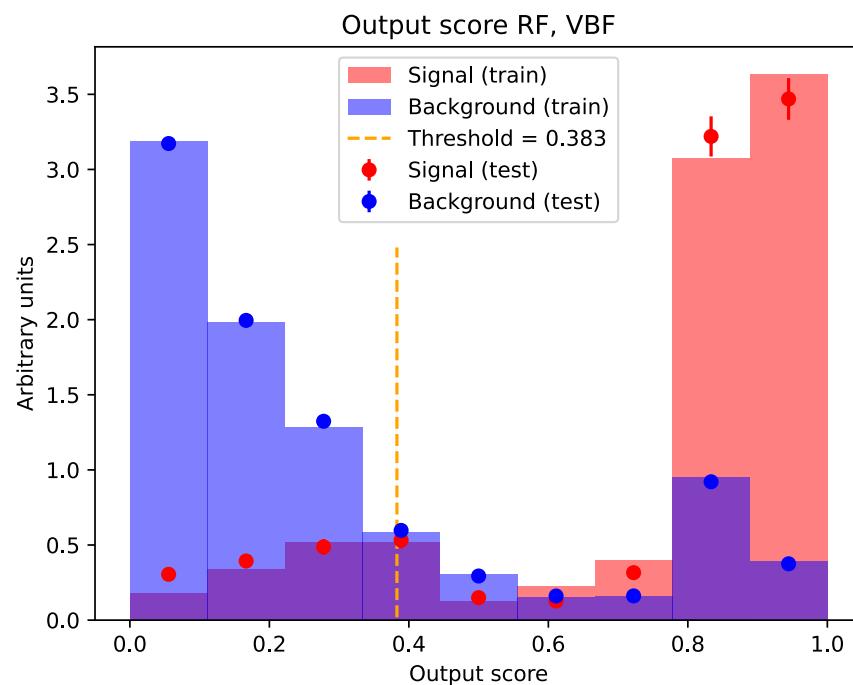
## Main function, P-E curves



```
1 # Get precision and recall
2 purity_test_rf, recall_test_rf, t_test_rf =
3 precision_recall_curve(y_true=y_test, probas_pred=y_prediction_test_prob[:, -1], sample_weight=W_test)
4
5 purity_train_rf, recall_train_rf, t_train_rf =
6 precision_recall_curve(y_true=y_train, probas_pred=y_prediction_train_prob[:, -1], sample_weight=W_train)
7
8 index_f_RF = fscore_threshold(purity_test_rf, recall_test_rf, t_test_rf)
9
10 # Threshold selection according to signal channel
11 cut_RF = 0.
12 if CHOICE == "ggH":
13     cut_RF = thresholds_test_rf[index_g_RF]
14 else:
15     cut_RF = t_test_rf[index_f_RF]
16
17 plotting_purity_vs_efficiency(f'RF_{CHOICE}', t_test_rf, t_train_rf, recall_test_rf,
18                                recall_train_rf, purity_test_rf, purity_train_rf, index_f_RF)
19 plotting_output_score(f'RF_{CHOICE}', y_train[:, 0], y_prediction_train_prob[:, 1],
20                      y_test[:, 0], y_prediction_test_prob[:, 1], cut_RF)
```

# RF Machine Learning: classificazione segnale/fondo (4)

## Output score

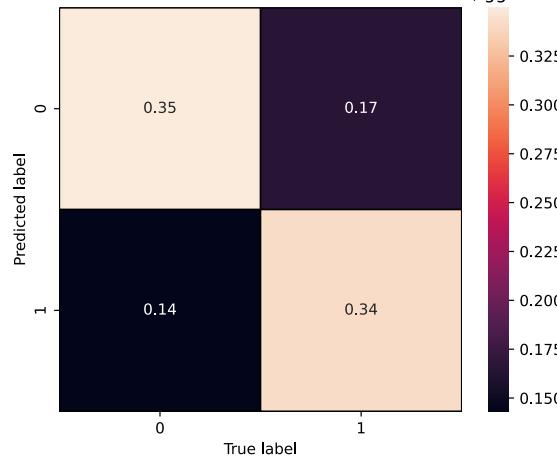


# RF Machine Learning: classificazione segnale/fondo (5)

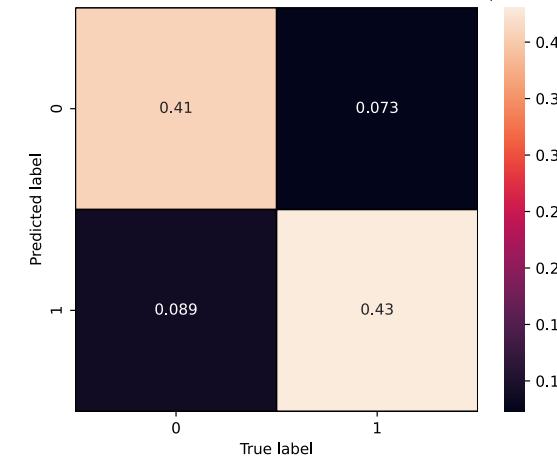
## Main function, results and confusion matrix

```
1 # Transform predictions into an array of 0, 1
2 filter_test_sig_rf = y_prediction_test_prob[:, -1] >= cut_RF # classify as signal
3 filter_test_bkg_rf = y_prediction_test_prob[:, -1] < cut_RF # classify as background
4 y_prediction_test_rf[filter_test_sig_rf] = 1
5 y_prediction_test_rf[filter_test_bkg_rf] = 0
6
7 # Metrics values for the ANN algorithm having fixed an ANN score threshold
8 accuracy_rf = accuracy_score(y_test[:, 0], y_prediction_test_rf, sample_weight=W_test[:, 0])
9 precision_rf = precision_score(y_test[:, 0], y_prediction_test_rf, sample_weight=W_test[:, 0])
10 recall_rf = recall_score(y_test[:, 0], y_prediction_test_rf, sample_weight=W_test[:, 0])
11 print(f'Threshold on the RF output : {cut_RF:.3f}')
12 print(f'RF Test Accuracy: {accuracy_rf:.3f}')
13 print(f'RF Test Precision/Purity: {precision_rf:.3f}')
14 print(f'RF Test Sensitivity/Recall/TPR/Signal Efficiency: {recall_rf:.3f}')
15 plotting_confusion_matrix(f'RF_{CHOICE}', y_test[:, 0], y_prediction_test_rf, W_test[:, 0])
```

Normalized Confusion Matrix for the test dataset RF, ggH



Normalized Confusion Matrix for the test dataset RF, VBF



...  
[Parallel(n\_jobs=4)]: Done 500 out of 500 | elapsed: 3.2s finished  
  
Accuracy of the RF model: 0.724 (test dataset)  
Accuracy of the RF model: 0.728 (train dataset)  
Best Threshold=0.456, G-Mean=0.691  
Best Threshold=0.143, F-Score simplify=0.535  
Threshold on the RF output : 0.456  
RF Test Accuracy: 0.691  
RF Test Precision/Purity: 0.704  
RF Test Sensitivity/Recall/TPR/Signal Efficiency: 0.673

...  
[Parallel(n\_jobs=4)]: Done 500 out of 500 | elapsed: 2.9s finished  
  
Accuracy of the RF model: 0.806 (test dataset)  
Accuracy of the RF model: 0.804 (train dataset)  
Best Threshold=0.397, G-Mean=0.838  
Best Threshold=0.383, F-Score simplify=0.708  
Threshold on the RF output : 0.383  
RF Test Accuracy: 0.838  
RF Test Precision/Purity: 0.829  
RF Test Sensitivity/Recall/TPR/Signal Efficiency: 0.855

[ML\\_Higgs.py](#),

[confusion plots](#)

# Machine Learning: classificazione segnale/fondo

## Main function & plotting function, variables plots

```
1 X_test = SC.inverse_transform(X_test)
2
3 # Plotting of ML variables with models predictions
4 for s in ML_dict[CHOICE]['ML_VARS']:
5     index_var = ML_dict[CHOICE]['ML_VARS'].index(s)
6     plotting_physical_variables(CHOICE, s, X_test[:, index_var], y_test[:, 0], y_prediction_test,
7                                 y_prediction_test_rf)
```

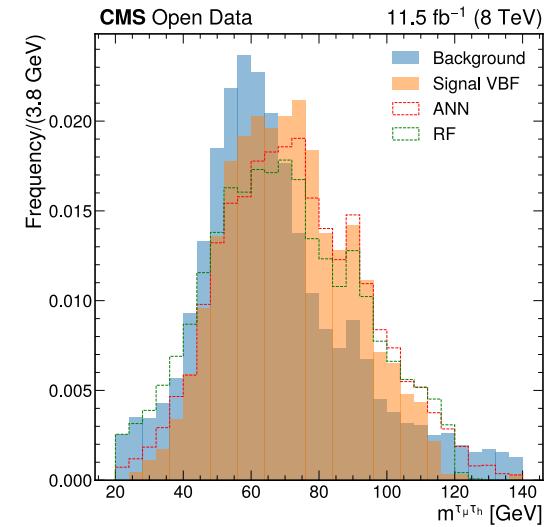
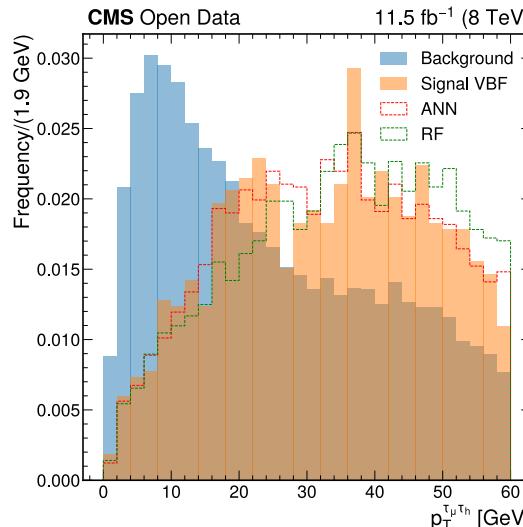
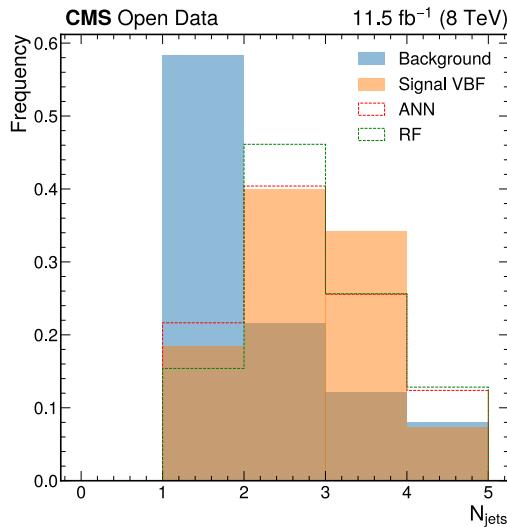
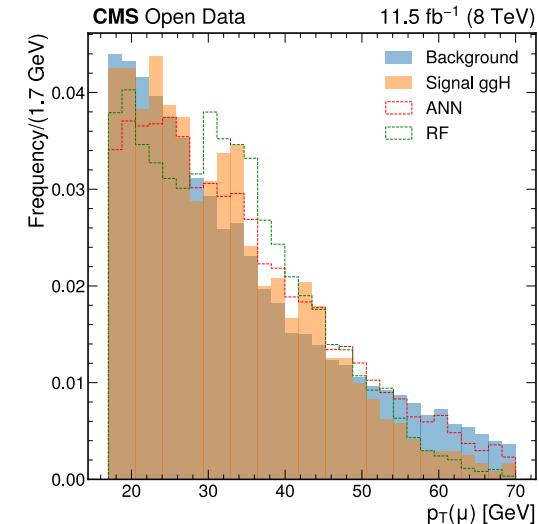
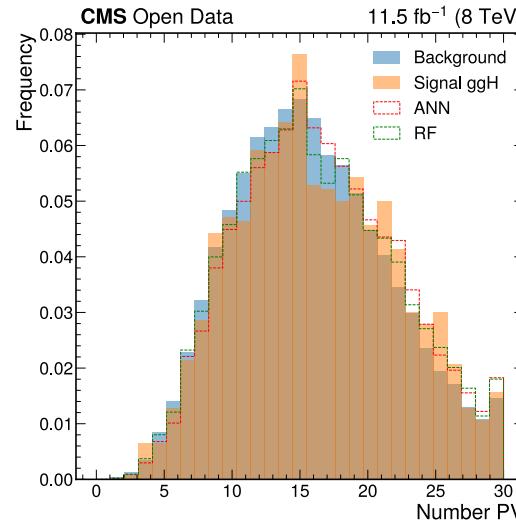
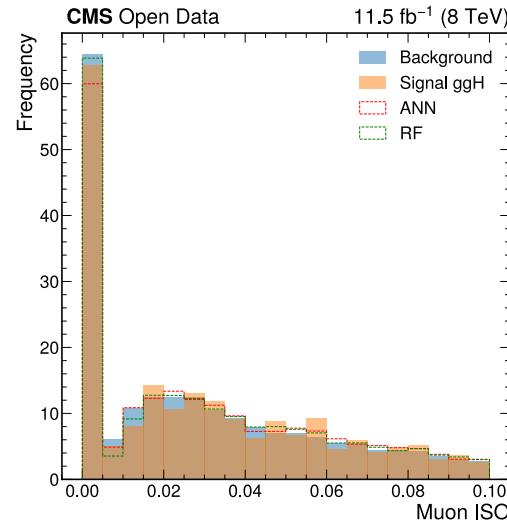
```
1 def plotting_physical_variables(CHOICE, s, feature, y_test, y_prediction_test, y_prediction_test_rf):
2     """ doc """
3     plt.style.use(hep.style.CMS) # or ATLAS/LHCb2
4     fig = plt.figure()
5     plt.hist(feature[y_test == 0], label='Background', density=1,
6             bins=ML_infoplot.plot[s]['bins'], alpha=0.5, linestyle='--')
7     plt.hist(feature[y_test == 1], label=f'Signal {CHOICE}', density=1,
8             bins=ML_infoplot.plot[s]['bins'], alpha=0.5, linestyle='--')
9     plt.hist(feature[y_prediction_test[:, 0] == 1], histtype='step', label='ANN', density=1,
10            bins=ML_infoplot.plot[s]['bins'], color='red', linestyle='--')
11    plt.hist(feature[y_prediction_test_rf == 1], histtype='step', label='RF', density=1,
12            bins=ML_infoplot.plot[s]['bins'], color='green', linestyle='--')
13    plt.ylabel(ML_infoplot.plot[s]['ylabel'])
14    plt.xlabel(ML_infoplot.plot[s]['xlabel'])
15    plt.legend()
16    plt.text(0.01, 1.02, r'$\bf CMS$ Open Data', transform=plt.gca().transAxes)
17    plt.text(0.65, 1.02, r'11.5 $fb^{-1}$ (8 TeV)', transform=plt.gca().transAxes)
18    plt.savefig(f'{NAME_DIR}/{s}_{CHOICE}.pdf', format='pdf')
19    plt.close(fig)
```

[ML\\_Higgs.py](#),

[plot\\_functions.py](#),

[Documentation](#)

# Machine Learning: classificazione segnale/fondo



Variables plots

# Machine Learning: Unittest

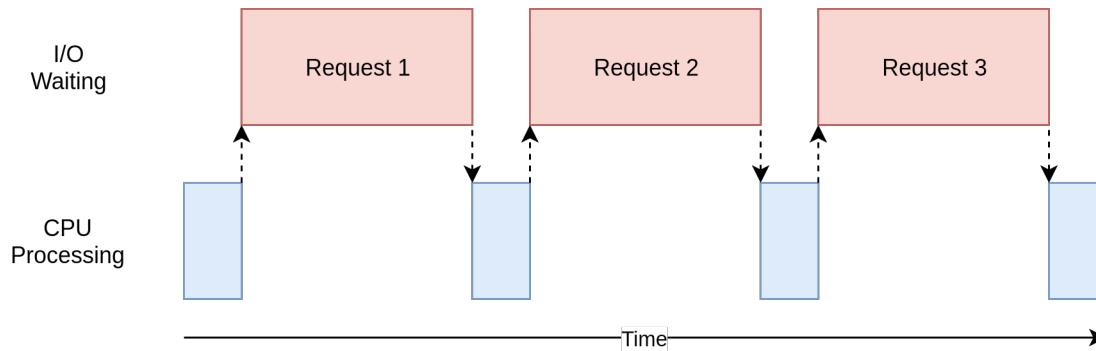
```
1 import unittest
2 import os
3 from ML_Higgs import read_root
4
5 class TestMLfiles(unittest.TestCase):
6     def setUp(self):
7         self.root_files_list = ["GluGluToHToTauTau", "VBF_HToTauTau", "DYJetsToLL",
8                                "TTbar", "W1JetsToLNu", "W2JetsToLNu", "W3JetsToLNu"]
9
10    self.signal_ggH = 4492
11    self.signal_VBF = 5326
12    self.bkg_TTbar = 13630
13
14    self.plots_list = ["nGoodJets_VBF", "PV_npvs_ggH", "MET_phi_ggH", "pt_vis_VBF",
15                       "muon_pt_VBF", "muon_eta_ggH", "muon_phi_VBF", "muon_iso_VBF",
16                       "tau_pt_ggH", "tau_eta_VBF", "tau_phi_ggH"]
17
18    def test_files_exist(self):
19        """ doc """
20        print("TEST 1 - Called...")
21        PATH = "ROOT_workspace/"
22        results = []
23        for file in self.root_files_list:
24            results.append(os.path.isfile(PATH + file + "_selected.root"))
25        self.assertTrue(all(results), "There aren't all ROOT files")
26
27    def test_read_root_ggH(self):
28        """ doc """
29        print("TEST 2 - Called...")
30        dataframe = read_root("GluGluToHToTauTau")
31        value_counts = dataframe['event'].value_counts().to_numpy()
32        self.assertEqual(value_counts[1], self.signal_ggH, "Must be 4492")
```

```
● ● ●
% python ML_test.py -b
test_files_exist (__main__.TestMLfiles)
Check on the existence of the file comparing with those in the root files list. ... ok
test_plots (__main__.TestMLfiles)
Check on the names of the saved plots comparing with those in plots list. ... ok
test_read_root_TTbar (__main__.TestMLfiles)
Check on the number of background events, e.g. in TTbar channel. ... ok
test_read_root_VBF (__main__.TestMLfiles)
Check on the number of signal events in the VBF channel. ... ok
test_read_root_ggH (__main__.TestMLfiles)
Check on the number of signal events in the ggH channel. ... ok
-----
Ran 5 tests in 13.381s
OK
```

```
32    def test_read_root_VBF(self):
33        """ doc """
34        print("TEST 3 - Called...")
35        dataframe = read_root("VBF_HToTauTau")
36        value_counts = dataframe['event'].value_counts().to_numpy()
37        self.assertEqual(value_counts[1], self.signal_VBF, "Must be 5326")
38
39    def test_read_root_TTbar(self):
40        """ doc """
41        print("TEST 4 - Called...")
42        dataframe = read_root("TTbar")
43        value_counts = dataframe['event'].value_counts().to_numpy()
44        self.assertEqual(value_counts[1], self.bkg_TTbar, "Must be 13630")
45
46    def test_plots(self):
47        """ doc """
48        print("TEST 5 - Called...")
49        for file in self.plots_list:
50            self.assertTrue(os.path.exists(f"ML_plots/{file}.pdf"),
51                           f"{file} not found!")
52
53 if __name__ == '__main__':
54     unittest.main(verbosity=2)
```

# Download dataset

## Sequential



```
% python download_sequential.py
The VBF_HToTauTau.root download is assigned to the function call: 1
The 'download' function runtime is: 28.9870s
The GluGluToHToTauTau.root download is assigned to the function call: 2
The 'download' function runtime is: 22.4422s
The 'sequential' function runtime is: 51.4292s
Downloaded files: ['VBF_HToTauTau.root', 'GluGluToHToTauTau.root']
```

```
1 # import libraries
2 import time
3 import requests
4 import functools
5
6 def timer(func):
7     """ doc """
8     @functools.wraps(func)
9     def wrapper_timer(*args, **kwargs):
10         start_time = time.perf_counter()
11         value = func(*args, **kwargs)
12         end_time = time.perf_counter()
13         run_time = end_time - start_time
14         print(f"The '{func.__name__}' function runtime is: {run_time:.4f}s")
15         return value
16
17     return wrapper_timer
18
19 @timer
20 def download(url, call):
21     """ doc """
22     response = requests.get(url, stream=True)
23     file_name = url.split('/')[-1]
24     print(f"The {file_name} download is assigned to the function call: {call}")
25     with open(file_name, "wb") as file:
26         for data in response.iter_content(1024):
27             file.write(data)
28     return file_name
29
30 @timer
31 def sequential():
32     """ doc """
33     urls = [
34         "https://root.cern/files/HiggsTauTauReduced/VBF_HToTauTau.root",
35         "https://root.cern/files/HiggsTauTauReduced/GluGluToHToTauTau.root"
36     ]
37     files = []
38     call = 0
39     for url in urls:
40         call += 1
41         files.append(download(url, call))
42     return files
43
44
45 if __name__ == "__main__":
46     done = sequential()
47     print(f"Downloaded files: {done}")
```

# Download dataset

## Parallel & Threading

```

1 import multiprocessing
2 ...
3 def parallel():
4     """ doc """
5     urls = [
6         "https://root.cern/files/HiggsTauTauReduced/VBF_HToTauTau.root",
7         "https://root.cern/files/HiggsTauTauReduced/GluGluToHToTauTau.root"
8     ]
9     start_time = time.time()
10    print(f"ID parent process: {os.getpid()}")
11    with multiprocessing.Pool() as pool:
12        pool.map(download, urls)
13    run_time = time.time() - start_time
14    print(f"Downloaded {len(urls)} in {run_time:.4f} seconds")
15    return "Done"

```

```

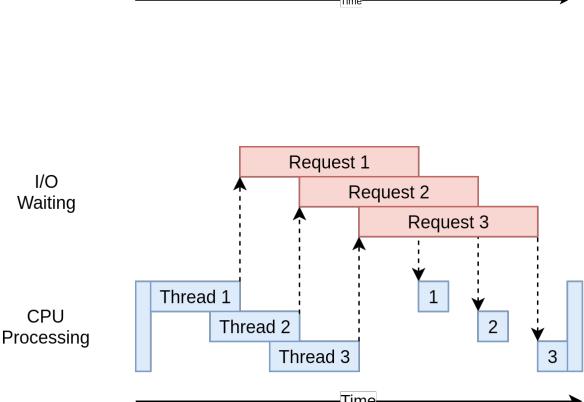
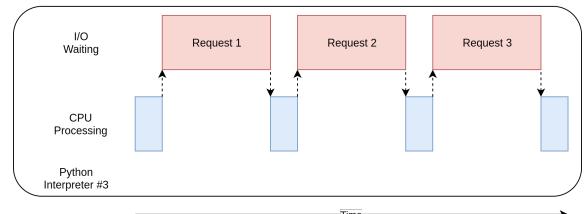
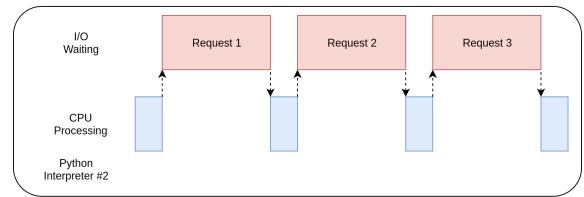
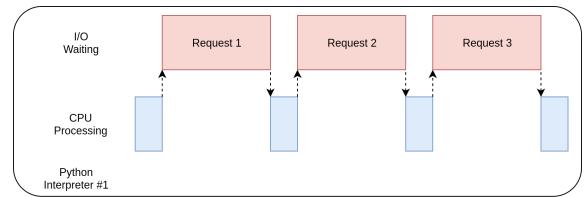
...
% python download_parallel.py
ID parent process: 2949
The VBF_HToTauTau.root download is assigned to the process with ID: 2951
The GluGluToHToTauTau.root download is assigned to the process with ID: 2953
The 'download' function runtime is: 30.7368s
The 'download' function runtime is: 34.8277s
Downloaded 2 in 35.5345 seconds
Done

```

```

...
% python download_threads.py
The VBF_HToTauTau.root download is assigned to the thread: Thread-1
The GluGluToHToTauTau.root download is assigned to the thread: Thread-2
The 'download' function (Thread-2) start time is: 0.5794s, end time is: 37.7494
The 'download' function (Thread-1) start time is: 0.5730s, end time is: 40.1196
Downloaded 2 in 39.5469 seconds
Done

```

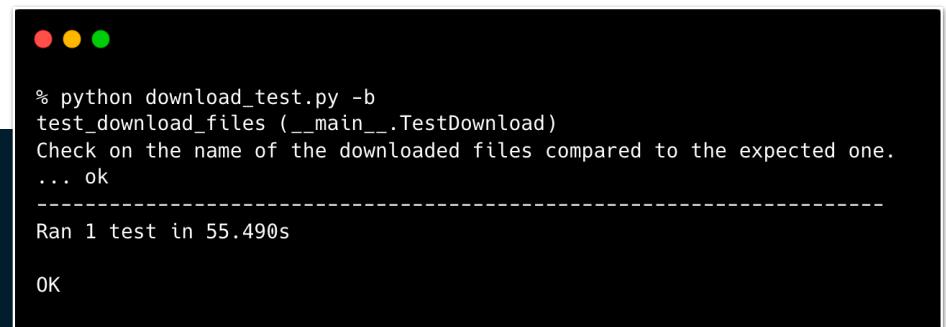


[download\\_parallel.py](#), [download\\_threads.py](#)

# Download dataset

## Unittest

```
1 import os
2 import unittest
3 from download_sequential import download
4
5 class TestDownload(unittest.TestCase):
6     def setUp(self):
7         self.url = "https://root.cern/files/HiggsTauTauReduced/"
8         self.files_list = ["VBF_HToTauTau.root", "GluGluHToTauTau.root"]
9         self.files_dim = [24184554, 20460281] # bytes
10
11    def test_download_files(self):
12        """ doc """
13        print("TEST 1 - Called...")
14        result = []
15        dim_result = []
16        call = 0
17        for file in self.files_list:
18            call += 1
19            result.append(download(self.url + file, call))
20            dim_result.append(os.stat(file).st_size)
21        self.assertEqual(result, self.files_list, "The files aren't the same")
22        self.assertEqual(dim_result, self.files_dim, "The dimensions aren't the same")
23
24 if __name__ == '__main__':
25     unittest.main(verbosity=2)
```



```
% python download_test.py -b
test_download_files (__main__.TestDownload)
Check on the name of the downloaded files compared to the expected one.
... ok
-----
Ran 1 test in 55.490s

OK
```

# Documentation - Sphinx website

<https://domric98.github.io/Searches-for-the-Higgs-boson-in-the-tau-tau-decay-channel/index.html>

Files prodotti per la documentazione:

- conf.py
- index.rst
- ML\_Higgs.rst
- plot\_function.rst
- Infoplot.rst
- ML\_test.rst
- download\_\*.rst

## Searches for the Higgs boson in the tau-tau decay channel

### Navigation

Machine Learning Contents:

- Code Documentation: [ML\\_Higgs.py](#)
- Code Documentation: [plot\\_functions.py](#)
- Code Documentation: [ML\\_infoplot.py](#)

# Welcome to Searches for the Higgs boson in the tau-tau decay channel's documentation!

Machine Learning Contents:

- Code Documentation: [ML\\_Higgs.py](#)
- Code Documentation: [plot\\_functions.py](#)
- Code Documentation: [ML\\_infoplot.py](#)
- Code Documentation: [ML\\_test.py](#)

Download Contents:

- Code Documentation: [download\\_sequential.py](#)
- Code Documentation: [download\\_parallel.py](#)
- Code Documentation: [download\\_threads.py](#)
- Code Documentation: [download\\_test.py](#)

# Workflow On GitHub

```
# This workflow will install Python dependencies, run tests and lint
# with a single version of Python
# For more information see: https://help.github.com/actions/language-
# and-framework-guides/using-python-with-github-actions

name: build

on:
  push:
    branches: [ main ]
  pull_request:
    branches: [ main ]

permissions:
  contents: read

jobs:
  build:
    runs-on: ubuntu-latest

    steps:
      - uses: actions/checkout@v3
      - name: Set up Python 3.10
        uses: actions/setup-python@v3
        with:
          python-version: "3.10"
      - name: Install dependencies
        run: |
          python -m pip install --upgrade pip
          pip install flake8 pytest
          if [ -f requirements.txt ]; then pip install -r requirements.txt; fi
      - name: Lint with flake8
        run: |
          # stop the build if there are Python syntax errors or undefined
          # names
          flake8 . --count --select=E9,F63,F7,F82 --show-source --
          statistics
          # exit-zero treats all errors as warnings. The GitHub editor is
          127 chars wide
          flake8 . --count --exit-zero --max-complexity=10 --max-line-
          length=127 --statistics
      - name: Test with pytest
        run: |
          pytest
```

The screenshot shows the GitHub Actions interface for the repository `DomRic98 / Searches-for-the-Higgs-boson-in-the-tau-tau-decay-channel`. The `Actions` tab is selected. Under the `Workflows` section, the `build` workflow is highlighted. Below it, the `python-app.yml` file is shown. The `36 workflow runs` table lists the following runs:

Event	Status	Branch	Actor
Update README.md	Success	main	Viooola
build	Success	main	DomRic98
Update README.md	Success	main	Viooola
Add plotting C++ code	Success	main	DomRic98
Add analysis C++ code	Success	main	DomRic98

File di configurazione del “build workflow”