

Measurement of the misidentification rate of electrons as hadronic taus in CMS experiment

Lorenzo Giannessi, Università degli studi di Firenze, Italy
DESY Summer Student Programme 2019
CMS experiment

Supervisor: Andrea Cardini

September 4, 2019

Abstract

Tau leptons produced during proton-proton collisions in LHC decay before reaching the inner layers of the CMS detector. Because of that they can only be reconstructed through their decay products. Particles can produce signatures very similar to the one produced by the τ decay products. An algorithm is thus used to match the particles reconstructed in CMS to a tau. On top of that other discriminators are needed to correct for the misidentification of other objects (jets, muon or electrons) as tau leptons. In this report an analysis concerning electrons mimicking hadronic τ is presented. The rate with which an electron is misidentified as a τ_h (i.e. the $e \rightarrow \tau_h$ fake rate, in which we call τ_h a τ that decays hadronically) is here measured with a "Tag & Probe" method. From the measurement a scale factor is obtained, which has to be applied on Montecarlo simulations of the $Z \rightarrow ee$ process in order for MC to model this misidentification rate and improve data/MC agreement.

Contents

1	Introduction	3
2	Motivation	3
3	Particle reconstruction and identification	4
3.1	Trigger	4
3.2	The particle flow (PF) algorithm	4
4	Hadronic τ reconstruction	5
4.1	The Standard Model - Brief review	5
4.2	The τ lepton	5
4.3	The Hadron Plus Strip algorithm	7
4.3.1	Electron faking τ_h	7
4.4	The DeepTau discriminator	8
4.4.1	Working points definition	8
5	The Tag & Probe method	9
5.1	Event preselection	9
5.1.1	Stitching and data driven methods	10
5.2	Pass/Fail regions	11
5.2.1	Data to MC disagreement	11
5.3	The Maximum Likelihood fit	13
5.3.1	Nuisance parameters	13
5.3.2	Binned Maximum Likelihood	15
5.3.3	Combine	17
6	Results	17
7	The (old) MVA discriminator	19
8	Acknowledgements	21

1 Introduction

Among all the particles that can be generated in a proton-proton collision, just a few of them reach the detector before decaying. All the others need to be reconstructed using algorithms that combine energy deposits and tracks to identify the particle. It is a priority to know how much these algorithms are reliable and what is the probability for these algorithms to fail (i.e. to misidentify a particle). Moreover, it is extremely important, for any kind of HEP analysis, to have Montecarlo simulations that accurately reproduce data, and to do so these simulations have to model also some potential "mistakes" in particle reconstruction. It is therefore necessary to measure the rates with which particles are misidentified, also in order to find scale factor to be applied on MC simulations for a better MC-to-data agreement.

2 Motivation

By getting a glimpse of the Higgs boson possible decay modes (as predicted by the Standard Model) it should be at least curious that the channels in which it was first discovered by the CMS and ATLAS collaborations were $H \rightarrow \gamma\gamma$ and $H \rightarrow ZZ \rightarrow 4l$ (4 lepton)[2, 3, 4], given that their branching ratios are not so large compared to other decay modes.

Higgs decays at $m_H=125\text{GeV}$

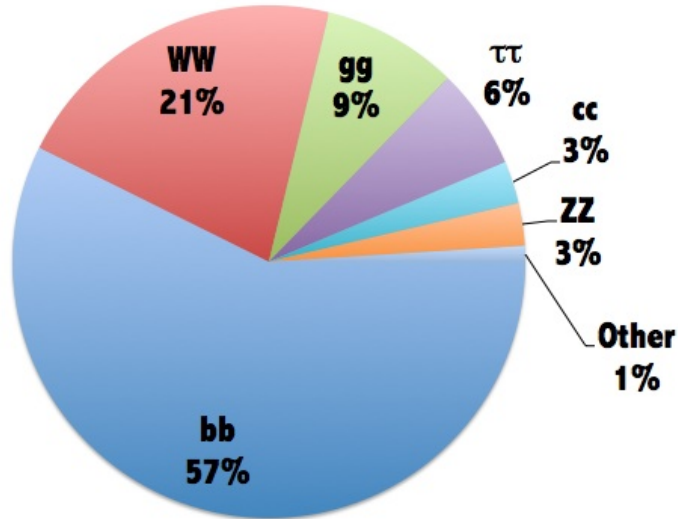


Figure 1: The Higgs boson decay channels and relative branching ratios. $H \rightarrow \gamma\gamma$ belongs to the "Other 1%" fraction, and the predicted branching ratio is 0.2%

However those two channels present very clear signatures compared to the ones with higher BRs, which in turn require more complex methods to reject the dominating backgrounds. It is now interesting and useful to study the other decay modes of the Higgs boson in order to check SM predictions. An important process to be studied is $H \rightarrow \tau\tau$, in the first place to search for deviations from the SM (BSM physics), as well as to study the coupling of the Higgs boson to leptons, given that it is the most sensitive channel among all the leptonic Higgs decay modes. For the analysis of $H \rightarrow \tau\tau$ various channels are studied, depending on the decay products of the taus:

- μe : fully leptonic
- $\mu\tau_h$ and $e\tau_h$: semi-leptonic
- $\tau_h\tau_h$: fully hadronic
- $\mu\mu$ and ee : di-muon and di-electron not studied because of high background

The measurement presented in this report is done in the $e\tau_h$ channel, in which the main backgrounds are $Z \rightarrow \tau\tau$ (as an irreducible background), $Z \rightarrow ee$, which plays an important role because of the misidentification of electrons as hadronic taus, and processes like W+Jets or QCD, in which a jet is reconstructed as an hadronic tau.

To study this Higgs decay mode in the $e\tau_h$ channel what we need in the first place is an algorithm for the reconstruction of τ leptons.

3 Particle reconstruction and identification

As a starting point for the identification of particles in CMS we need to know which signature such particles should leave in the detector. The signature is based on the physical properties of the particle and the detector response. Moreover, for such particles as the τ , which decays before reaching the first layer of detectors, it is even more difficult to find an algorithm that recognizes them properly. We need to have an efficient method to discriminate between different particles signatures.

3.1 Trigger

Due to the huge amount of events that take place in CMS, only a small fraction of them, which is interesting to analyse, are saved; this first selection is made by the CMS trigger system[7]. This is done in order to reduce the amount of data that have to be stored, and to discard events which are uninteresting for a specific study.

3.2 The particle flow (PF) algorithm

After an event has been selected and saved, it has to be processed: the raw information (i.e. tracker hits, calorimeter deposits etc.) have to be combined to produce higher

level objects on which data analysis can be more easily performed. This is done by the Particle-Flow (PF) algorithm[1], which uses the information taken from the trackers and calorimeters to extract the physics objects which will be used in the analysis (we can call them "PF candidates"). Once an object has been identified, the corresponding tracks and clusters are removed from the subsequent searches.

The PF algorithm is used to reconstruct particles which interact directly with the detectors, but our analysis concerns the τ leptons, which obviously don't belong to this category. In order to reconstruct them one has to consider all the possible final states that can be produced by this particle and then look for the PF candidates that match these final states. More in details, this analysis takes place in the $e\tau_h$ channel, therefore we expect a τ decaying leptonically¹ and the other decaying hadronically. In this study we are specifically interested in the reconstruction of the latter; we must take into account hadronic decay modes of the τ lepton and what signature they leave in CMS.

4 Hadronic τ reconstruction

We understood that the problem of τ reconstruction is a challenging one, we must start by presenting the properties of this particle and how we can use them to detect it.

4.1 The Standard Model - Brief review

The theory that describes all known elementary particles and forces is the Standard Model[5]. Since the τ lepton is predicted by this theory, is useful to provide a brief introduction on it. In this theory matter is made out of 12 elementary fermionic particles divided in two main groups: 6 quarks and 6 Leptons. These fermions are further grouped in 3 generations with increasing mass. Furthermore, there are 4 vector bosons and 1 scalar boson (the Higgs) are present; they are the mediators of the four fundamental interactions that this theory describes, while the Higgs boson is responsible for the electroweak symmetry breaking[6]. For what concerns leptons, there are 3 charged leptons: e , μ and τ , and 3 neutral leptons: ν_e , ν_μ and ν_τ .

4.2 The τ lepton

The τ lepton is a third generation charged lepton, it is the most massive with a mass of $1776.86 \pm 0.12 MeV$, and a mean life of approximately $2.9 \times 10^{-13}s$. It is the only one massive enough to decay hadronically, and it does so approximately 2/3 of the times. For every decay a tau neutrino is present in the final state, so we must expect some missing transverse energy in each process that include the production (and consequent decay) of a τ . Feynmann diagram in Fig.3 is useful to visualize τ decays. The decay always happens through weak interaction, specifically with a charged-current interaction (W^\pm are involved). To be more specific, τ decays leptonically in muons or electrons with

¹ $\tau \rightarrow e\nu_\tau\bar{\nu}_e$

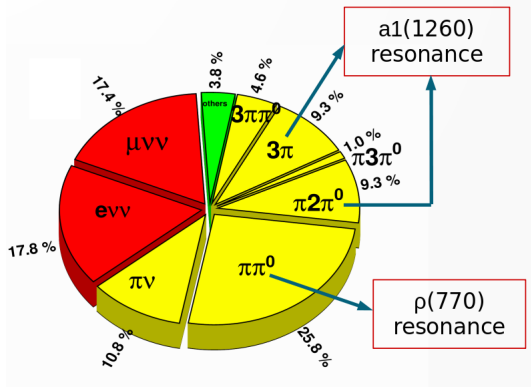


Figure 2: τ lepton decay modes and branching ratios

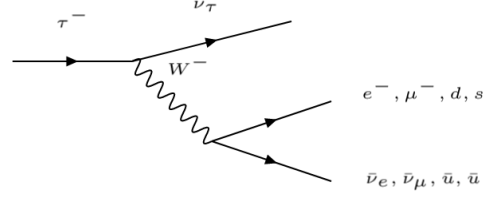


Figure 3: Feynmann diagrams of tau decay processes, notice the hadronic final state with π^- ($d\bar{u}$) and the rarer K^- ($s\bar{u}$)

35.2% branching ratio (17.4% muonic, 17.8% electronic), and 64.8% hadronically. The hadronic decays consist of combinations of charged and neutral mesons plus a ν_τ . Mesons involved in these decays are mostly pions and K s, nearly only π^0 for what concerns the neutral ones; a_1 and ρ resonances can take place. In Tab.1 all decay modes are shown in details together with their branching ratios and possible resonances. In Fig.2 leptonic modes and hadronic modes with pions in final state. The branching ratios are the same for the charge-conjugate processes.

Table 1: τ lepton decay modes and branching ratios

Decay mode	Resonance	Branching ratio (%)
Leptonic		35.2
$\tau^- \rightarrow e^- \bar{\nu}_e \nu_\tau$		17.8
$\tau^- \rightarrow \mu^- \bar{\nu}_\mu \nu_\tau$		17.4
Hadronic		64.8
$\tau^- \rightarrow h^- \nu_\tau$		11.5
$\tau^- \rightarrow h^- \pi^0 \nu_\tau$	$\rho(770)$	25.9
$\tau^- \rightarrow h^- \pi^0 \pi^0 \nu_\tau$	$a_1(1260)$	9.5
$\tau^- \rightarrow h^- h^+ h^- \nu_\tau$	$a_1(1260)$	9.8
$\tau^- \rightarrow h^- h^+ h^- \pi^0 \nu_\tau$		4.8
Other		3.3

4.3 The Hadron-Plus-Strip algorithm

The algorithm used to reconstruct hadronic taus is called Hadron Plus Strip (HPS). The decay products of the τ_h are mainly charged and neutral hadrons, so this algorithm should look for the signatures of these particles. A charged hadron leaves as a signature a track in the tracker and a deposit in ECAL and HCAL; it is also referred to as a *prong*. For what concerns the neutral hadrons they are π^0 nearly in 100% of the times, and are reconstructed using *strips*. A *strip* is a narrow region in the ECAL in which deposits are found, characterized by short width in the η direction but a large length in the ϕ direction. The π^0 produced from the τ decays with a mean life of $(8.52 \pm 0.18) \times 10^{-17} s$ in two photons²; at least one of the photons does pair production, the e^+e^- pair interacting with the magnetic field splits along the ϕ direction, and a strip-like object is therefore generated in the calorimeter. The algorithm looks for charged hadrons among all the PF candidates and for strips; then a matching with one of the relevant τ decay modes is required to reconstruct a τ_h . Moreover, the invariant mass of all charged hadrons and reconstructed π^0 s must be compatible with ρ or a_1 mass, for the decay channel in which there is a resonance. 4 different categories matching decay modes can be identified:

- **one prong** (h^\pm): the visible mass of the reconstructed τ_h should be compatible with the mass of a charged hadron (typically pion or K meson).
- **one prong plus one strip** ($h^\pm\pi^0$): the invariant mass should match m_ρ (770 MeV).
- **one prong plus two strips** ($h^\pm\pi^0\pi^0$): the invariant mass should match m_{a_1} (1260 MeV).
- **3 prongs** ($h^\pm h^\mp h^\pm$): this can match with two different decay modes: $\tau^- \rightarrow h^- h^+ h^- \nu_\tau$ and $\tau^- \rightarrow h^- h^+ h^- \pi^0 \nu_\tau$, depending on whether the invariant mass of the three prongs matches with m_{a_1} or not.

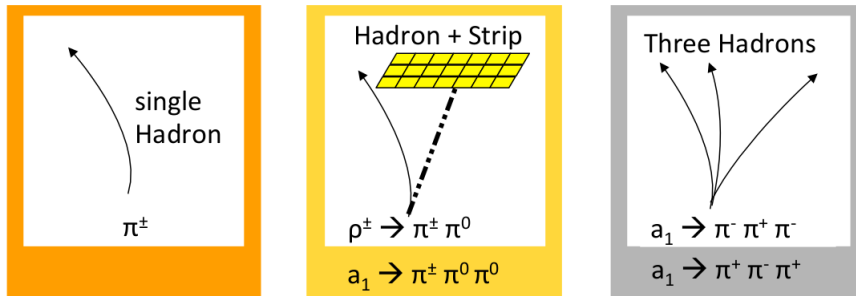


Figure 4: Some decay modes reconstruction of HPS algorithm

²To be more precise: the Branching ratio for $\pi^0 \rightarrow \gamma\gamma$ is $\sim 99\%$, and the $\sim 1\%$ Br of $\pi^0 \rightarrow e^+e^-\gamma$ is here neglected. Furthermore, the same effects happen with this decay products

It can incidentally happen that another particle matches one of these requirements and is reconstructed as a τ_h . The main source of contamination for the HPS algorithm are jets, muons and electrons, which are the subject of this analysis.

4.3.1 Electron faking τ_h

It can happen that an electron mimicks the signature of a charged hadron. Moreover, bremsstrahlung radiation can potentially take place while the electron rapidly slows down, and emitted photons can mimick a strip. For that reason, mainly for decay channels with one prong and one prong plus one strip, an electron can be labeled as an hadronic τ by the HPS algorithm.

In this case we need a discriminator to distinguish between true taus and electrons.

4.4 The DeepTau discriminator

It shouldn't be surprising that the discrimination from different types of objects that fake the τ_h signature is dependant on many variables, so a multivariate analysis is a natural choice for solving the problem. As soon as we have to deal with 3 sources of contamination, a Deep Neural Network with 4 possible outputs is used. The output neurons are: true τ_h , electron faking τ_h , muon faking τ_h and jet faking τ_h . With this setup three discriminators can be defined, depending on what source of contamination we want to distinguish the τ_h from³ (i.e.: against jet, against electron, against muon).

4.4.1 Working points definition

For each discriminant what we obtain is a test statistic in which signal (true τ_h) and background (one out of the three possible fakes) should be enough separated (for this analysis: DeepTau2017v2VSe). Based on a certain efficiency that we require, a cut is applied on the discriminant, defining two orthogonal regions. From that cut a *working point* is obtained. The larger is the efficiency that we want to have, the smaller is the background rejection, so the working point is looser; viceversa: the larger the background rejection, the tighter is the WP. Available working points are reported for the against electron discriminant, which is relevant for this analysis, with relative efficiencies (Tab.2).

Table 2: Working points and relative efficiencies

WP	VVTight	VTight	Tight	Medium	Loose	VLoose	VVLoose
Efficiency	60%	70%	80%	90%	95%	98%	99%

³e.g.: the anti-electron discriminator is calculated as $\frac{p_\tau}{p_\tau + p_e}$, in which p_τ stands for the true τ_h output, and so on.

5 The Tag & Probe method

As said in Section 4.3.1, there is a nonzero probability for an electron to leave a signature that can mimick the one left by an hadronic τ , the *Tag & Probe* method is suitable for measuring this fake rate.

5.1 Event preselection

The starting point for this analysis consists of the **EGamma** dataset for 2018, together with Montecarlo simulations of all relevant processes that can contribute to the final state of interest (*templates*):

- $Z \rightarrow ee$
- DY others: mainly $Z \rightarrow \tau\tau$
- Diboson
- W+Jets
- $t\bar{t}$

The QCD background is evaluated with the data driven ABCD method (Section 5.1.1).

What we first need is two samples of events, one coming from data, the other coming from Montecarlo simulations, in which we require to have two particle defined as following:

- The Tag: a well identified and isolated electron
- The Probe: an hadronic τ which has passed loose preselection criteria

A first event selection is operated with a *Single Electron Trigger*. We use the *Electron35* trigger (HLT_Ele35_WPTight_Gsf_v) that requires each event to have a well isolated electron with a transverse momentum of more than 35 GeV. Now, in addition to the Tag, which is the electron that fired the trigger, we require a particle that can be possibly matched with at least one among all the possible hadronic decay modes of the τ that the HPS algorithm can identify (our Probe)⁴. Further cuts are then applied on the samples:

- on the Tag
 - i) Electron $P_T > 35\text{GeV}$
 - ii) Electron $|\eta| < 2.1$
 - iii) The electron has to match with the particle which fired the trigger in a cone of $\Delta R < 0.5$

⁴From now on the words "electron" and "tag" will be synonyms, as well as "hadronic tau" and "probe"

- iv) Electron relative isolation $I_e < 0.1$; I_e is a variable used to discriminate isolated electrons from electrons produced in jets
- on the Probe
 - i) hadronic tau $P_T > 20 GeV$
 - ii) hadronic tau $|\eta| < 2.3$
 - iii) The τ_h has to pass the DeepTau discriminator against Jets (Tight) and Muons (Loose) before being checked against the electron. These WP choice is based on working region commonly used in $H \rightarrow \tau\tau$ analysis
 - iv) VVVLoose DeepTau discriminator against electron
- transverse mass $m_T < 30 GeV$, in order to avoid most of W+Jets background
- $\Delta R > 0.5$ between tag and probe, in order to not have two collimated particles

By defining the probe particle this way we imply that it has passed "loose" preselection criteria

5.1.1 Stitching and data driven methods

Some procedures to improve the MC-to-data agreement are performed:

For what concerns the simulated processes in which different number of jets can take place (i.e. W+Jets) the stitching procedure has to be done. For this processes several MC samples are available: the ones generated with a particular number of jets (e.g.: W+1Jet, W+2Jets, W+3Jets and W+4Jets), and the "inclusive" samples, in which we have all the different cases together. The stitching procedure consists of extracting from the inclusive sample smaller subsets with a defined number of jets, by matching the number of jets at generator level, then adding each subset with the independently generated sample with the same number of jets and weighting each sample with the cross section of that process. The aim of this procedure is to increase the statistic as well as improving Montecarlo accuracy, as it is easier to obtain better precision for what concerns detector performances in processes with a defined number of jets.

To estimate the QCD background a data driven ABCD method is used: This background is evaluated as $data - AllMCs$ in a control region ("anti-iso" region), then a scale factor between the "same sign, anti-iso" and the "opposite sign, anti-iso" regions is calculated. Given that the two variables (isolation and sign) are not correlated, the same scale factor can be used to calculate the QCD background in the "opposite sign, iso" region, which is the region in which the measurement is done, given the background in the "same sign, iso" region.

At this point our samples are finally ready to undergo the selection.

5.2 Pass/Fail regions

For both Montecarlo simulations and data, the sample is now divided in two orthogonal regions, which are **PASS** and **FAIL**, depending on whether the probe particle has passed or not the anti-electron discriminant which we are testing. The number of events populating these two regions will depend on where the cut was applied on the `tauagainstEleRawDeepTau`, i.e. on the working point. For each working point we have a boolean variable in the NTuple which is true if the probe has passed the anti-electron discriminant, false otherwise.

As an example visible mass plots in the pass region are shown (Fig.5).

Notice how increasing the working point's tightness is reducing the relative contribution of $Z \rightarrow ee$, until the Z boson mass peak vanishes. It is also interesting to notice the displacement of the Z boson mass peak for $Z \rightarrow ee$ and $Z \rightarrow \tau\tau$. In the former we don't have neutrinos, while in the latter both taus decay, one leptonically ($\tau \rightarrow e\nu_\tau$), the other one hadronically, generating two neutrinos; we therefore expect missing energy, so that the yellow peak is located at lower (visible) mass.

If a MC simulated $Z \rightarrow ee$ event passes the discriminant, it means that one electron has been misidentified as a τ_h . It means that we can just calculate the *misidentification rate* as the fraction of electrons which pass the anti-electron discriminator:

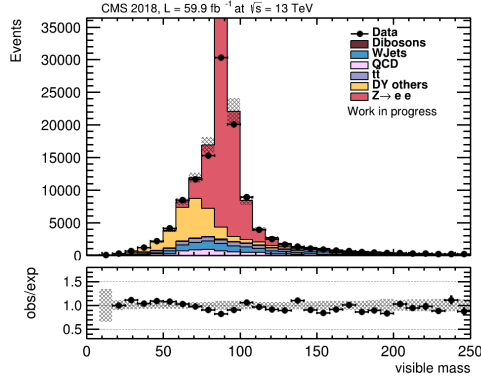
$$\epsilon = \frac{N_{Z \rightarrow ee}^{pass}}{N_{Z \rightarrow ee}^{pass} + N_{Z \rightarrow ee}^{fail}} \quad (1)$$

This is the *pre-fit fake rate*.

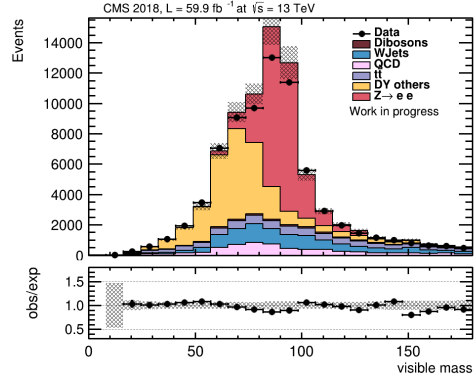
Remark: $Z \rightarrow ee$ is the process on which the pre-fit measurement is done since it's basically the only one where there are electrons faking taus in the $e\tau$ channel.

5.2.1 Data to MC disagreement

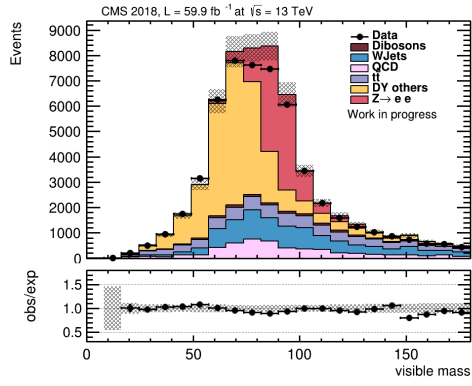
When both MC and data undergo the anti-electron discriminant we would expect that the number of events in the pass and in the fail region should be the same for MC and data, within statistical uncertainties. By looking at the plots for the pass region, we can notice that the agreement between data and MC needs to be improved. To improve it we should measure the "*real*" fake rate on data and then apply the correct scale factor on the process of interest. In fact this MC-to-data disagreement is partially or completely removed when the scale factor r is measured and applied. r represent the ratio between the *post-fit*(ϵ') and the *pre-fit*(ϵ) fake rate.



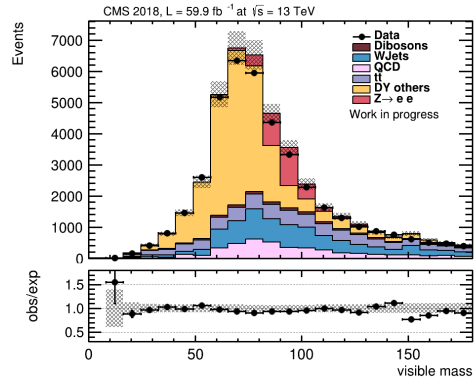
(a) VLoose



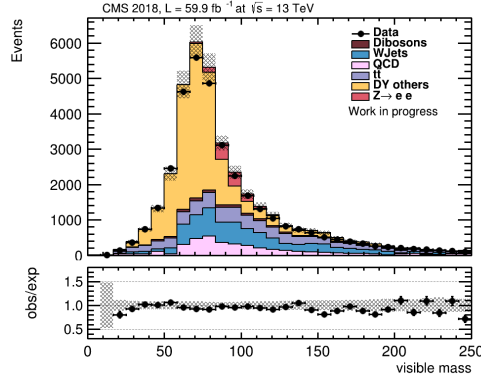
(b) Loose



(c) Medium



(d) Tight



(e) VTight

Figure 5: Visible mass plots for various WPs. Both MC simulated processes and data are plotted

It will then be applied as a weight on simulated $Z \rightarrow ee$ events in order to account for this discrepancy. In Fig.6 plots for many variables in the pass region for the Medium WP are shown: notice that the disagreement between data and MC often overcomes the uncertainty bands.

5.3 The Maximum Likelihood fit

Until now we didn't use data to measure the fake rates, since Eq.1 only uses MC simulations of signal events, and represent how MC samples behave with the anti-electron discriminator. The "real" misID rate has to be measured on data, and is the *post-fit* fake rate. We want to obtain a measurement of the scale factor $r = \frac{\epsilon'}{\epsilon}$, so that we are able to calculate the *post-fit* fake rate, and we manage to do that using a multi-parameter maximum likelihood fit in which the parameters can be labeled as:

- The parameter of interest (POI) $r = \frac{\epsilon'}{\epsilon}$
- An array of nuisance parameters θ which are used to account for systematic uncertainties in the MC simulations.

Pre-fit values for all the parameters are given to the fitting program as inputs, then, moving in the parameters space a maximum in the likelihood function is found, so that the post-fit values are obtained. The variable chosen for the fit in this analysis is the reconstructed invariant mass of the electron and all visible decay products of the τ_h (*visible mass*). In order to account for differences in performances and geometry of the calorimeters, two regions of $|\eta|$ are defined: one relative to the Barrel calorimeter: $|\eta| < 1.460$ and the other relative to the Endcap: $|\eta| > 1.558$. A narrow gap of $|\eta|$ is left uncovered because the performance of the detector abruptly drops in that region.

5.3.1 Nuisance parameters

A way to take into account systematic uncertainties that affects our samples is to introduce them as nuisance parameters for the fit. Nuisance parameter are other arguments of the Likelihood function that are not POI, but influence the likelihood model. In Tab.3 & 4 a list of the nuisance parameters which are used in this analysis is shown, with the relative template on which are applied and the starting value for the fit. We can in general recognise two different types of systematic uncertainties: *shape* uncertainties, which can for instance modify the position or the height of a peak in a certain process, and *normalization* uncertainties, which instead are scale factor for the total number of events for a simulated process.

Notice (Fig.7) that, as the fit is performed, not only $Z \rightarrow ee$ is modified, but also other background processes are. That happens because all processes are affected by systematic uncertainties, and if a systematic uncertainty for a certain process is modified, also the number of events in each bin for that process changes.

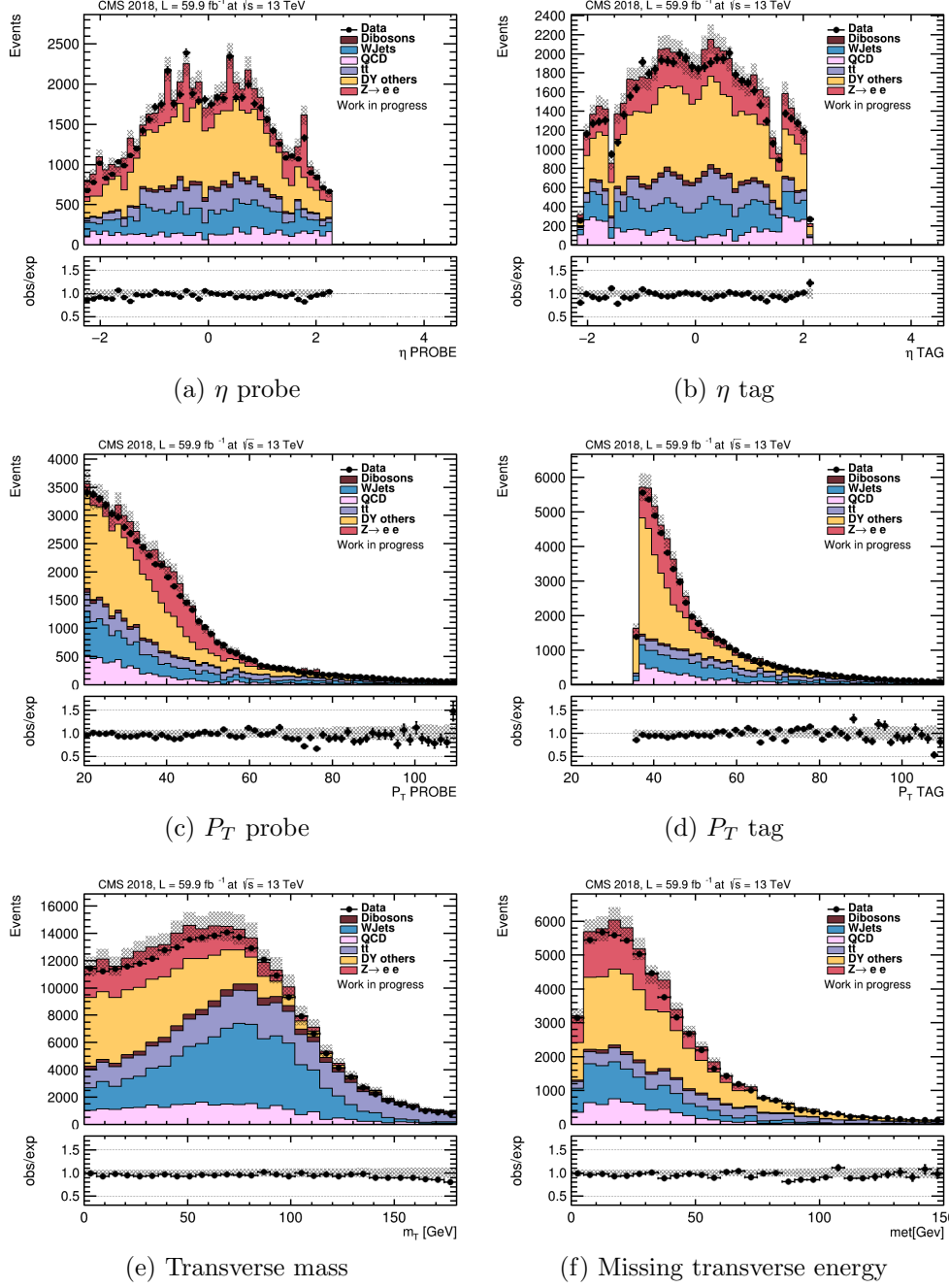


Figure 6: Plots of several variables for data and MC samples. $\frac{data}{MC}$ shown in the bottom plots. Shaded bands represent the statistical error on montecarlo samples due to finite number of events, together with systematic uncertainties, while vertical bars (barely visible) represent the statistical uncertainty on data[♣]

♣ Transverse mass definition: $m_T = \sqrt{P_T^{ele} E_T^{miss} (1 - \Delta\Phi_{ele,met})}$

Table 3: Shape uncertainties

Uncertainty	Affected processes	Pre-fit value
Electron energy scale	$Z \rightarrow ee$	1% B, 2.5% E ⁵
τ_h energy scale	$Z \rightarrow \tau\tau$	1.5%
$e \rightarrow \tau_h$ energy scale	$Z \rightarrow ee$	3%
Visible mass resolution	$Z \rightarrow ee$	20%

Table 4: Normalization uncertainties

Uncertainty	Affected processes	Pre-fit value
Integrated luminosity	All MCs	2.6%
Electron isolation/identification/trigger	All MCs	2%
Tau identification	All MCs	3%
$t\bar{t}$ cross section	$t\bar{t}$	10%
Diboson and single-top cross section	Diboson	10%
W+Jets normalization	W+Jets	20%
QCD normalization	QCD	20%
DY normalization	DY Others	3%
$Z \rightarrow ee$ normalization	$Z \rightarrow ee$	6%

5.3.2 Binned Maximum Likelihood

In this particular case we are not trying to fit a histogram (MCs) on a smooth function which depends on some parameters, but on another histogram (data) instead. What we need is a Likelihood function as shown in Eq. 2

$$\mathcal{L}(r, \boldsymbol{\theta}) = \prod_{i=1}^{N_{bins}} \text{Poisson}(n_i | \nu(r, \boldsymbol{\theta})) \times p(\tilde{\boldsymbol{\theta}} | \boldsymbol{\theta}) \quad (2)$$

in which:

- i is the index running on bins of the histograms
- N_{bins} is the total number of bins
- $\tilde{\boldsymbol{\theta}}$ is the estimator of the parameter array $\boldsymbol{\theta}$

⁵E: Endcap, B: Barrel

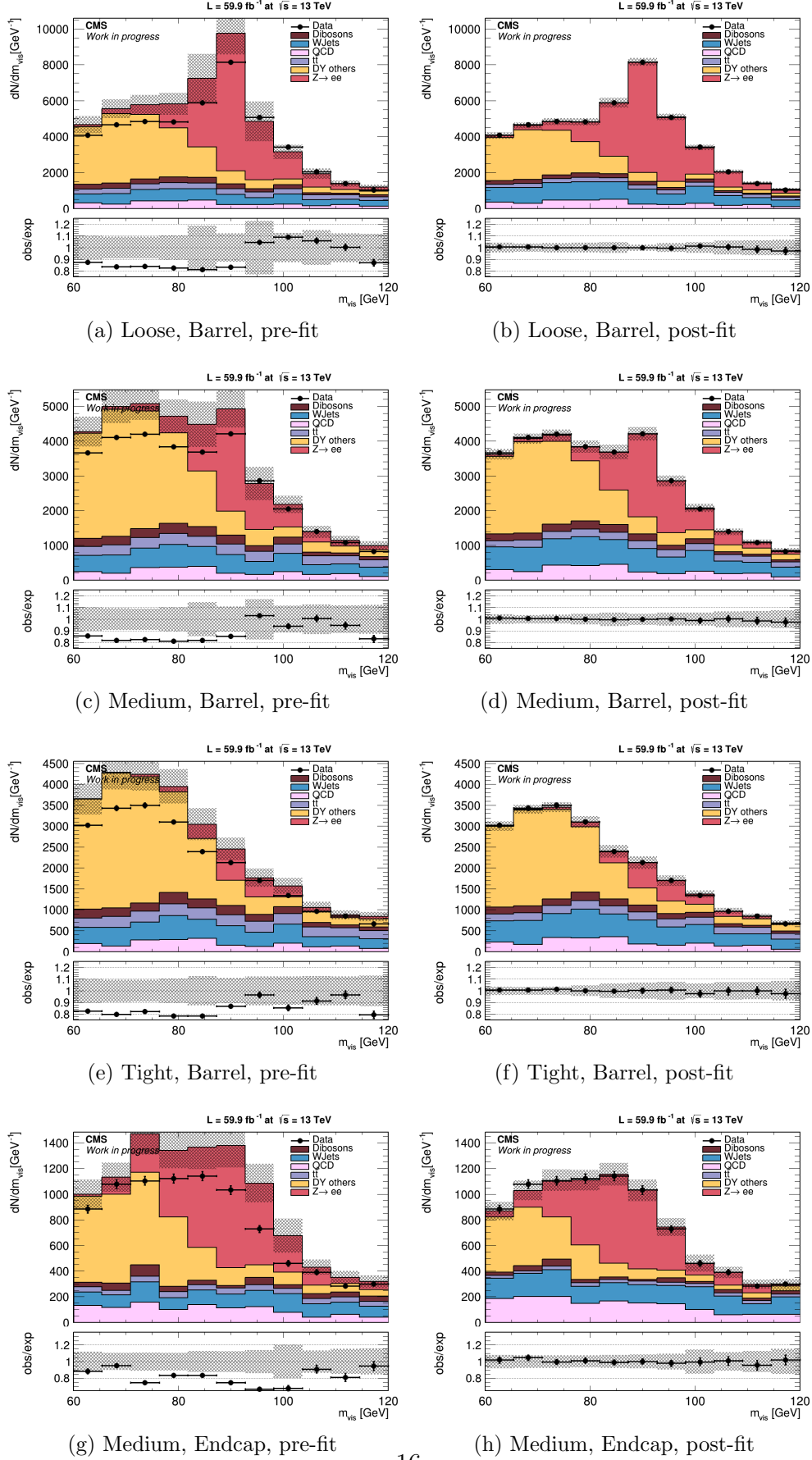


Figure 7: Pre-fit / Post-fit plots. Notice how all background processes vary and not only $Z \rightarrow ee$, especially in tighter WPs, in which the relative contribution of signal events is smaller

- $Poisson(n_i|\nu_i(r, \boldsymbol{\theta}))$ is the Poisson distribution with mean value $\nu_i(r, \boldsymbol{\theta})$, which is the number of events expected in the i -th bin by the Montecarlo simulation with that parameters.
- $p(\tilde{\boldsymbol{\theta}}|\boldsymbol{\theta})$ is the *a priori* distribution for the nuisance parameter array

Some plots are shown for various working points in Fig.7, in order to understand how the fit modifies MC samples.

5.3.3 Combine

COMBINE is the program which we use for the fit. It takes as inputs the starting points for all parameters, which are the values listed for the systematic uncertainties, and a pre-fit POI r , usually chosen close to 1. The program maximizes the likelihood function in the parameter space, with expedients to avoid stopping in relative maxima, and then gives back the parameters that maximize the function.

6 Results

The output of the fitting program is the scale factor r , together with the post-fit values for all nuisance parameters. In tables 5 and 6 fake rates measurements are shown.

Table 5: Scale factors and fake rates in Barrel region. Error bars for post-fit FR have been symmetrized

$$postfitFR = prefitFR \times scalefactor$$

BARREL ($ \eta < 1.460$)	pre-fit FR	post-fit FR	scale factor
VVLoose	$(4.5 \pm 0.5) \times 10^{-2}$	$(5.8 \pm 0.7) \times 10^{-2}$	$r = 1.280^{+0.014}_{-0.015}$
VLoose	$(2.4 \pm 0.3) \times 10^{-2}$	$(3.3 \pm 0.4) \times 10^{-2}$	$r = 1.368^{+0.029}_{-0.033}$
Loose	$(1.0 \pm 0.1) \times 10^{-2}$	$(1.31 \pm 0.18) \times 10^{-2}$	$r = 1.31^{+0.05}_{-0.05}$
Medium	$(3.9 \pm 0.4) \times 10^{-3}$	$(5.0 \pm 0.9) \times 10^{-3}$	$r = 1.35^{+0.09}_{-0.09}$
Tight	$(1.14 \pm 0.13) \times 10^{-3}$	$(1.5 \pm 0.4) \times 10^{-3}$	$r = 1.32^{+0.22}_{-0.22}$
VTight	$(4.8 \pm 0.5) \times 10^{-4}$	$(7 \pm 3) \times 10^{-4}$	$r = 1.4^{+0.4}_{-0.4}$
VVTight	$(2 \pm 0.2) \times 10^{-4}$	$(2.1 \pm 0.7) \times 10^{-4}$	$r = 1.4^{+0.9}_{-0.9}$

These scale factor values are the first step for a CMS recommendation, to be applied on MC simulated $Z \rightarrow ee$ process for analysis using 2018 data that involve reconstructed τ_h objects.

For almost every working point and in both Barrel and Endcap, the measured scale factor exceeds one, which means that MC simulations were underestimating the fake rate before fit. By looking at the scale factor we can clearly see a trend in the uncertainty, which increases as the WP gets tighter. While the results for looser working points are

Table 6: Scale factors and fake rates in Barrel region. Error bars for post-fit FR have been symmetrized

$$postfitFR = prefitFR \times scalefactor$$

ENDCAP ($ \eta > 1.558$)	pre-fit FR	post-fit FR	scale factor
VVLoose	$(8.7 \pm 0.9) \times 10^{-2}$	$(1.14 \pm 0.13) \times 10^{-1}$	$r = 1.318^{+0.016}_{-0.017}$
VLoose	$(4.4 \pm 0.5) \times 10^{-2}$	$(5.8 \pm 0.8) \times 10^{-2}$	$r = 1.314^{+0.033}_{-0.039}$
Loose	$(1.9 \pm 0.2) \times 10^{-2}$	$(2.7 \pm 0.5) \times 10^{-2}$	$r = 1.38^{+0.07}_{-0.08}$
Medium	$(9 \pm 1) \times 10^{-3}$	$(1.21 \pm 0.12) \times 10^{-2}$	$r = 1.35^{+0.13}_{-0.12}$
Tight	$(2.7 \pm 0.3) \times 10^{-3}$	$(4 \pm 1) \times 10^{-3}$	$r = 1.51^{+0.27}_{-0.29}$
VTight	$(9.3 \pm 1.1) \times 10^{-4}$	$(6.5 \pm 6.5) \times 10^{-4}$	$r = 0.7^{+0.8}_{-0.7}$
VVTight	$(4.2 \pm 0.5) \times 10^{-4}$	$(4.2 \pm 4.2) \times 10^{-4}$	$r = 1.0^{+1.6}_{-1.0}$

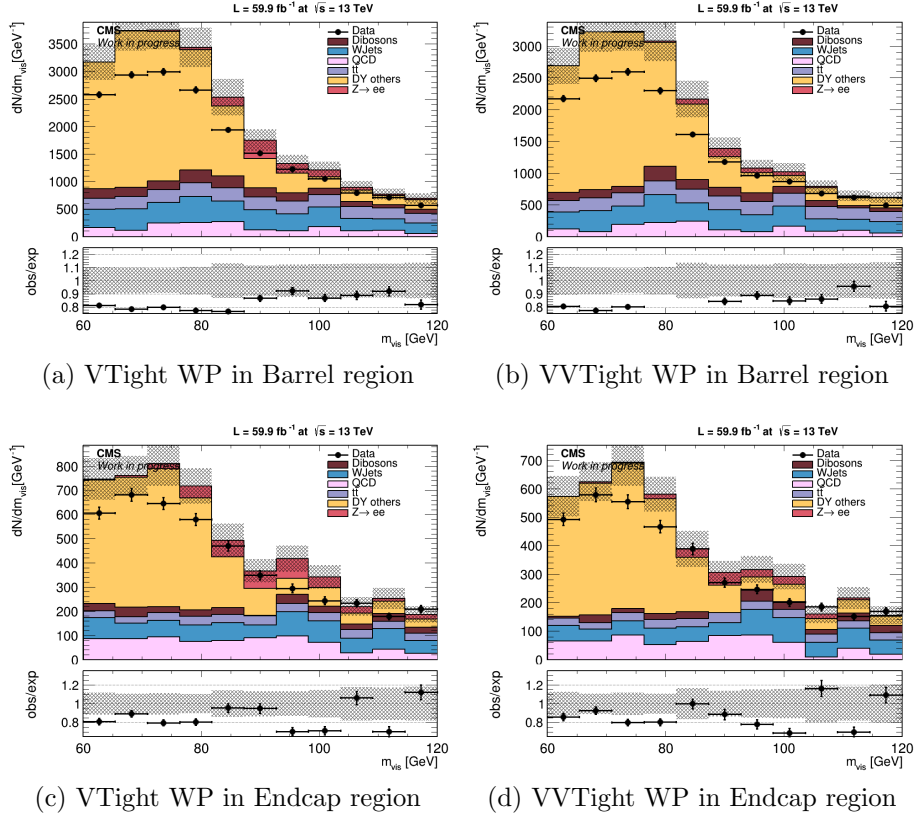


Figure 8: pre-fit plots for tightest working points

satisfactory⁶, if we look at the scale factors for VTight and VVTight WPs we see that the uncertainty is too large: in the barrel region the relative uncertainty reaches almost

⁶Relative uncertainties acceptably small

65%, in the endcap region the result is even consistent with zero. We can explain this type of result by looking at the pre-fit plots of the visible mass for VTight and VVTight working points (Fig.8).

As the working point become more and more tight the relative contribution of signal events rapidly decreases, as the background rejection of the anti-electron discriminator increases. At some point the amount of signal events is so little that the value of the POI r doesn't really matter for the fit, as the MC-to-data disagreement is mainly due to other process and not $Z \rightarrow ee$, which contribute with a negligible number of events. In general, we can suppose that the trend in uncertainties could be related to the decreasing total number of events, which makes statistical uncertainties larger (look at the y-axis scale in Fig.8c and Fig.8d).

To come to conclusion, a first look at the new DeepTau discriminator was given, fake rates and scale factor for $Z \rightarrow ee$ were measured, with quite accurate results for all WP below VTight. For the VTight and VVTight WPs we can conclude that the background rejection of the discriminator is so high that the number of electron that end up in the pass region is almost negligible.

7 The (old) MVA discriminator

In this study we gave a first look at the DeepTau discriminator in the $e\tau_h$ channel, measuring its fake rates and scale factors. The Tag & Probe was already successfully used to measure these scale factors for the previous τ_h discriminator used in CMS, that is referred to as the *MVA discriminator* (MVARun2v1DBoldDMwLT). It was also based on a multivariate analysis, but it used a Boosted Decision Tree instead of a DNN. It was interesting during this analysis to see how differently these two discriminators behave on the samples. In this small section we make a comparison between MVA and DeepTau discriminators for anti-electron and anti-jet.

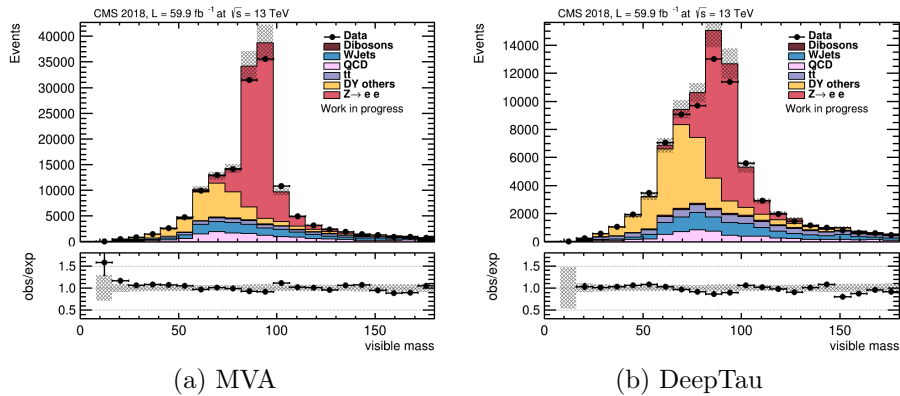


Figure 9: visible mass for Loose WP against electron, Tight WP against jet

At a first glance (Fig.9) we see that the background rejection of the DeepTau is clearly

greater than the MVA one for the same WP. Just by looking at the y-axis scale it is clear that the number of events in the pass region is strongly reduced for nearly every process. Especially for the $Z \rightarrow ee$ process, it looks like the ability to reject the background is improved.

Some interesting features were found in the $|\eta|$ plots (Fig.10), in which the shape is considerably changing between the two discriminators. This could be due to the different eta ranges used during the training of the neural network. Notice the absence of a net distinction between the Barrel and Endcap region which was instead clearly evident in the MVA.

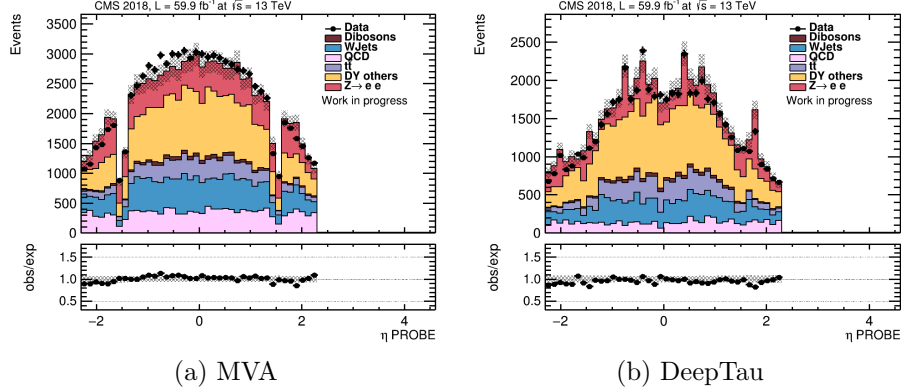


Figure 10: $|\eta|$ of the probe for Medium WP against electron, Tight WP against jet

We then noticed that, for what concerns the against jet discriminator, the Tight WP for MVA is more similar to the Medium for DeepTau than the Tight (Fig.11). This is another evidence that DeepTau is more "powerful" (in terms of background rejection) than MVA, for both against-electron and against-jet discriminators.

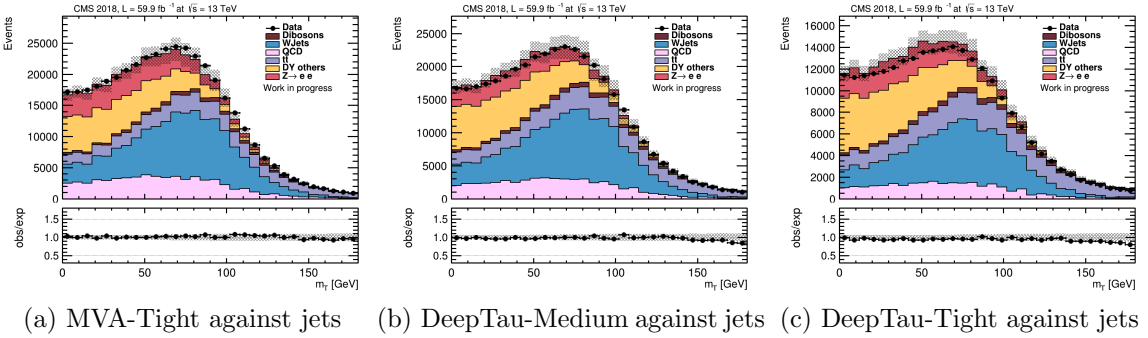


Figure 11: transverse mass, Tight WP against electron

8 Acknowledgements

I want to say a big thank you to all my colleagues from the Higgs-to-tau-tau group for their advices and support expressed in every monday morning meeting (which I enjoyed a lot), for their kindness and comprehension; to my office colleagues Valeria, Antonio and Ilya, to have made the work environment enjoyable and informal at any time. My supervisor Andrea deserves a special thank for being not only a professor who taught me a lot, but a nice friend too. Finally, thank you to all the summer students with whom I had super fun, and were really friendly.

References

- [1] A. Sirunyan et al., Particle-flow reconstruction and global event description with the CMS detector, *Journal of Instrumentation*, vol. 12, pp. P10003P10003, oct 2017.
- [2] CMS Collaboration. Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. *Phys. Lett. B* 716, 2012. doi:10.1016/j.physletb.2012.08.021. arXiv:1207.7235.
- [3] ATLAS Collaboration. Observation of a new particle in the search for the standard model Higgs boson with the ATLAS detector at the LHC. *Phys. Lett. B* 716, 2012. doi:10.1016/j.physletb.2012.08.020. arXiv:1207.7214.
- [4] CMS Collaboration. Observation of a new boson with mass near 125 GeV in pp collisions at $s = 7$ and 8 TeV. *JHEP* 06 (2013) 081, 2013. doi:10.1007/JHEP06(2013)081. arXiv:1303.4571.
- [5] Brian R. Martin, Graham Shaw; "Particle Physics, 4th Edition", WILEY, 2017
- [6] P. W. Higgs, Broken symmetries and the masses of gauge bosons, *Phys. Rev. Lett.* 13 (1964) 508, doi:10.1103/PhysRevLett.13.508.
- [7] CMS Collaboration. The CMS trigger system. doi: 10.1088/1748-0221/12/01/P01020. arXiv: 1609.02366
- [8] Yiwen Wen, "Tau Lepton: as a Tool for Hunting Standard Model and Beyond Standard Model Higgs", 2019