Studies on the Separation of Prompt and Tau Muons with the ATLAS Detector using Machine Learning Methods







Physics at the LHC Seminar, Daniel Ortmann, 7/25/2022







Outline

- Motivation]
- **Approach and Data** 2
- **Feature Variables and Classifier Architecture** 3
- **Training Set Bias in** *p***_T** 4
- **Training Strategies** 5
- **Generalization to Other Processes** \mathbf{b}
- **Feature Importance** 7
- **Application on BSM Physics Signal** 8
- **Conclusion and Outlook** 9



1 Motivation



- Process of interest: decay of the tau lepton τ into the muon μ (mediated by the W boson)
- Conservation of lepton family number \longrightarrow neutrinos in the final state



on μ (mediated by the *W* boson) the final state



1 Motivation

- Muons originating from tau decays $(W, Z \rightarrow \tau + ... \rightarrow \mu + ...)$, referred to as *tau muons* are reconstructed as prompt muons (W, $Z \rightarrow \mu + ...$)
- Tau decay length $l_{\tau} \approx 0.09 \text{mm} \longrightarrow \tau$ is not directly trackable by detector
- How can we distinguish between those ? •
- Closest approach of the muon track to the beam spot perpendicular to the beam axis: d_0 \bullet
- Current analyses use only d_0 to distinguish between tau and prompt muons \longrightarrow use other variables as well
- Potential application in several analyses:
 - Verification of the Lepton-Flavour Universality: $R(\tau/\mu) = B(W \to \tau \nu_{\tau}) / B(W \to \mu \nu_{\mu}) \approx 1$
 - Measure $B(W \to \tau \nu_{\tau} \to \mu \nu_{\mu} \nu_{\tau} \nu_{\tau})$ and $B(W \to \mu \nu_{\mu})$ using a dedicated classifier
 - Extrapolate $B(W \to \tau \nu_{\tau} \to \mu \nu_{\mu} \nu_{\tau} \nu_{\tau})$ to $B(W \to \tau \nu_{\tau})$ using the branching ratio of taus to muons $\approx 1/6$
 - Application in the search for BSM physics (supersymmetry in particular)



2 Approach and Data

- Approach: Supervised Learning with deep neural networks (DNNs)
- Datasets: Muons satisfying the *Medium* WP and originating from different Monte-Carlo simulated processes (with POWHEG-BOX interfaced with PYTHIA 8 as Event-Generator):
 - $t\bar{t} \rightarrow bb + qq' + \ell \nu_{\ell}$ and $bb + \ell \nu_{\ell} + \ell' \nu_{\ell'}$ $(\ell = \tau, \mu)$
 - $Z \to \ell \ell$
 - $W \to \ell \nu_{\ell}$
 - $W^* \to \tau \nu_{\tau} \longrightarrow$ in order to generate a large amount of high- p_{T} tau muons
- No isolation WP is used to suppress fake muons, since they are optimized for prompt but not tau muons
- Signal: tau muon μ_{τ}

Background: prompt muons µ and fake muons (objects misclassified as muons or non-prompt muons, e.g. hadron decay in jets) • No distinction between fake and prompt muons is made (tertiary classification was tested and no improvement was observed)





3 Feature Variables and Classifier Architecture

 $\rho' =$

Eight features were selected for their distinctiveness between tau muons and background



$$\Delta z_0 = |z_{\mathbf{pv}} - z_0|$$

momentum related	Others
p_{T}	$ \eta $
$rac{\Delta p_{\mathrm{T}}}{p_{\mathrm{T}}}$	Calorimeter Isolation
ho'	
$\frac{p_{\rm T}^{\rm ID} - p_{\rm T}^{\rm MS}}{p_{\rm T}^{\rm CB}}$	$\eta = -\ln \tan(\frac{\theta}{2})$



3 Feature Variables and Classifier Architecture

- DNN is implemented using *Keras* a library on top of *TensorFlow*
- Preprocessing: transform each input variable x_j on same scale (mean \overline{x}_j and standard deviation σ_j over all samples) $x_j \mapsto x'_j = \frac{x_j - \overline{x}_j}{\sigma_j}$
- 3 layers containing 128 (fully-connected) neurons
- Introduce non-linearities: ReLU(x) = max(0, x)
- Batch normalization after and Dropout on last neuron layer
- Sigmoid function $S(z) = \frac{1}{1 + e^{-z}} \longrightarrow$ squeeze values into [0, 1]
- Binary crossentropy loss $L(y, y^{\text{true}}) = -y^{\text{true}} \log(y) - (1 - y^{\text{true}}) \log(1 - y)$
- $t\overline{t}$ dataset is split up into
 - Train set: used for training
 - Test set: evaluated during training to monitor possible overfitting
 - Evaluation set: evaluated after training for analysing the performance of the trained classifier



sible overfitting e performance of the trained classifier



4 Training Set Bias in p_T

Most tau muons are at low p_{T} and the high- p_{T} regime is dominated by prompt muons \longrightarrow high- p_{T} muons are with high ulletprobability classified as prompt muons





4 Training Set Bias in p_T

- As an illustration, train a classifier without a dedicated handling of this bias ightarrow
- ightarrow
- Remaining $t\bar{t}$ samples are defined to be evaluation set ightarrow
- Use efficiency as performance measurement, i.e. the fraction of samples per class surviving an output score threshold ullet
- Use fixed signal efficiency of 50 % for better comparison ullet



Sample 250,000 tau muons and 125,000 prompt and fake muons from $t\bar{t}$, then split into train (75 %) and test dataset (25 %)

- Here the bias, i.e. signal efficiency drop at high p_{T} , is ightarrowequivalent of an increasing high- p_{T} prompt muon efficiency
- Small amount of high- p_T fake muons \longrightarrow classifier ignores them resulting in a high fake muon efficiency
- We need a training strategy that copes with that!
- Three strategies were used:
 - Flat sampling during batch calls
 - Sample weights
 - Distance correlation



5 Flat Sampling during Batch Calls

- Approach: Use a distribution flat in the entire p_{T} regime for the training set, to give high- p_{T} tau muons equally importance compared to background muons
- Flat sampling during batch calls: full availabe *tī* dataset is split into train and test set, where the latter is defined to be also the evaluation set
- During train time and at each epoch: a different subset is sampled from the train set that is flat in p_T and on which the DNN is trained on \longrightarrow improve generalization power
- 10,000 tau and 5,000 prompt as well as fake muons are sampled from 5 GeV p_T-bins at each epoch, last bin is an overflow bin [195,∞]
- If a subset does not contain enough muons of a particular category, muons may be selected multiple times



- Overall background is reduced:
 - At low $p_{\rm T}$ prompt muons rejection increased, fake muon rejection decreases only slightly
 - At high p_{T} fake muons are highly rejected at the expense of a slight increasement of the prompt muon efficiency







6 Generalization to Other Processes

• Output score thresholds are reused



• Significant less rejection power of fake muons from W/Z events than in $t\bar{t} \longrightarrow$ likely originates in different fake compositions



• Low process dependency \longrightarrow has seen full availabe $t\bar{t}$ dataset during training





6 Generalization to Other Processes



- $t\bar{t}$ tau muon efficiency is 0.5 (per definition)
- Signal efficiency drop up to 0.1 in the regime 10 GeV $\leq p_{T} \leq 60$ GeV. Small process dependency for larger p_{T}
- Strong W* efficiency drop for $p_T > 185 \text{ GeV} \longrightarrow W^*$ sample is dominated by high- p_T muons and the regime $p_T > 185 \text{ GeV}$ high- p_{T} regime

was not covered by enough muons during training which results in a strong efficiency drop \rightarrow train dedicated classifier for





7 Feature Importance

- Shuffle the *j*-th feature variable across all samples and keep the others fixed \rightarrow break its correlation with the other variables
- importance to the model
- Metric: area under ROC curve (AUC)
- Measurement: AUC unshuffeld $AUC_i^{shuffled}$ (normalized by the mean of all eight feature importances)



- ρ' could be dropped from the training without performance loss

• If model performance does not decrease substantially w.r.t. a certain metric \rightarrow shuffled feature variable provides only slight

• Fake muons are much less isolated than tau muons, hence the classifier relies heavily on the calorimeter isolation variable • Except of the calorimeter isolation, the other variables do not contribute much to the discrimination power compared to $|d_0|$



8 Application on BSM Physics Signal

- At least an extension of the SM is necessary to address big open questions in elementary particle physics \rightarrow supersymmetry (SUSY)
 - Could solve the hierarchy problem and would offer an Dark Matter candidate

 - Consider the Minimal Supersymmetric Standard Model (MSSM)
 - Some superpartners:
 - Tau lepton $\tau \longrightarrow stau \tilde{\tau}$
 - Neutral gauge bosons $W^0, B^0 \longrightarrow$ neutral gauginos \tilde{W}^0, \tilde{B}^0
 - Neutral components of the Higgs field $(H_u^0 \text{ and } H_d^0) \longrightarrow$ neutral higgsinos $(\tilde{H}_u^0 \text{ and } \tilde{H}_d^0)$

 - $\tilde{\chi}_1^0$ would be the Dark Matter candidate (LSP and R-parity)

- Hypothesized process of interest:
 - Stau pair production: $\tilde{\tau}\tilde{\tau} \to \tau\tau + \tilde{\chi}_1^0 \tilde{\chi}_1^0 \to \mu\nu_\mu\nu_\tau + \tau_{had} + \tilde{\chi}_1^0 \tilde{\chi}_1^0$
- Application of classifier for the search in the semileptonic decay channel:
 - Suppress $W + \text{jets} \rightarrow \mu \nu_{\mu} + \text{jets}$ (if jet is misidentified as τ_{had}

• Postulates the existence of partner particle for each particle in the SM (superpartner), which differs in spin by half a unit

• Electroweak symmetry breaking $\longrightarrow \tilde{H}^0_u, \tilde{H}^0_d$ and \tilde{W}^0, \tilde{B}^0 mix to four neutral mass eigenstates: neutralinos ($\tilde{\chi}^0_i, j=1, 2, 3, 4$)





8 Application on BSM Physics Signal

- $\Delta m := m_{\tilde{\tau}} m_{\tilde{\chi}_1^0}$
- In addition, evaluated on stau muons corresponding to a large variety of mass parameters ($m_{\tilde{\tau}} \in [80, 440]$ GeV and $m_{\tilde{\chi}_1^0} \in [1, 200] \,\mathrm{GeV})$
- Reuse output score thresholds (i.e. such that tagging efficiency of $t\bar{t}$ tau muons is 0.5)



• Evaluate on muons originating from stau-pair production (stau muons) in dependency of stau mass $m_{\tilde{\tau}}$ and mass splitting



8 Application on BSM Physics Signal

• BSM signal efficiency around 50% and independent of the $m_{\tilde{\tau}}$ and $m_{\tilde{\chi}_1^0}$ up to the first order





9 Conclusion and Outlook

- Tagging of muons originating from tau decays interesting for SM measurements as well as BSM searches
- Presented development of ML-based tagger trained on e.g. $t\bar{t}$ events \bullet
- Flat sampling during batch calls classifier shows best performance:

 - Evaluating the classifier for processes other than $t\bar{t}$: signal efficiency reduces up to 10 % \longrightarrow include other processes in the training as well



• At a p_{T} -flat signal efficiency of 50 %: prompt and fake muon efficiency of 12.5 % to 22.5 % and 5 % and 10 %, respectively

17

9 Conclusion and Outlook

- Before the application in an actual analysis could take place, several additional steps would be required first:

 - optimized for prompt muons not tau muons)
 - techniques in the electron channel \longrightarrow increase the sensitivity reach of a search for new physics

• ATLAS event simulation not perfect \rightarrow calibrate taggers, i.e. match the efficiency in simulation with the one measured in data • Isolation WPs on leptons are typically used in ATLAS analyses \longrightarrow study interplay between WPs and classifiers (WPs are

• Decays from taus to muons and electrons occur at the same rate \rightarrow develop classifier dedicated to electrons and apply same

18

*p*_T Distribution of the Simulated Processes

Sample Weights

- Same dataset as in the naive approach with the biased training set
- Determine sample weights in 5 GeV steps up to 60 GeV such that the distribution is flat in p_T

$$\sum_{m} total \longmapsto L_{m}' total = \frac{1}{m} \sum_{k=1}^{m} w_{k} \cdot L(y_{k}, y_{k}^{true})$$

and thus fluctuations in the loss surface (may spoil its convergence)

• Sample weights corresponding to the interval [55,60) GeV are applied on every sample with $p_T > 60$ GeV, to prevent large weights

- Same datasets used as before. no modification of the dataset is made, but rather on the loss function itself
- if output score and p_{T} of muon are not correlated, this may also ensure a similar tagging performance over the whole \bullet p_{T} regime \longrightarrow Add new term to the loss function
- Distance correlation $dCorr_n^2(X, Y)$ as measurement of correlation between two variables
- Modify loss function: $L_m^{\text{total}} \mapsto L_m^{\text{total}} = L_m^{\text{total}}(y, y^{\text{true}}) + \lambda \cdot d\text{Corr}_m^2(y, p_T)$, with tunable hyperparameter λ \bullet

- Fake muon rejection is improved but on the price of a lower prompt muon rejection
- Training set is still dominated by low- p_{T} muons and hence the decorrelation occurs mostly at low $p_{\rm T}$ as well

Observed n paired samples $(X, Y) = (x_i, y_i)_{i=1}^n$ with $a_{ij} = |x_i - x_j|$ and $b_{ij} = |y_i - y_j|$ define

$$A_{ij} = a_{ij} - \frac{1}{n} \sum_{i=1}^{n} a_{ij} - \frac{1}{n} \sum_{j=1}^{n} a_{ij} + B_{ij} = b_{ij} - \frac{1}{n} \sum_{i=1}^{n} b_{ij} - \frac{1}{n} \sum_{j=1}^{n} b_{ij} + \frac{1}{n}$$

The empirical distance covariance $dCov_n(X, Y)$ is then defined by

$$dCov_n^2(X, Y) = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij}$$

Finally, the discrete distance correlation $dCorr_n(X, Y)$ is given by

 $dCorr_n^2(X, Y) = \frac{dCov_n^2(X, Y)}{\sqrt{dCov_n^2(X, X)dCov_n^2(Y, Y)}}$

(b) $\lambda = 1$

(d) $\lambda = 50$

High-p_T Muon Taggers

• Train two high- p_T tau muon tagger exclusively for $p_T \ge 60 \text{ GeV} \longrightarrow$ increase importance for high- p_T muons

• Trained on same weighted dataset as the sample weight classifier, but exclusively for the regime $p_T > 60 \text{ GeV}$

• Trained on tau muons from W* and on prompt muons from $t\bar{t}$, fake muons are not considered

High-p_T Muon Taggers

• High- p_T muon tagger trained on tau muons from $t\bar{t}$ events improves performance only w.r.t. tau muons with $p_T > 285$ GeV

• sample weights classifier

• High- p_{T} muon tagger

27

High-p_T Muon Taggers

- High- p_{T} muon tagger trained on tau muons from W^* events
- tau muons provided by calorimeter isolation
- Bad generalization to tau muons from other processes

• Inappropriate for fake muon rejection \longrightarrow not trained on fake muons, i.e. could not learn distinctiveness between fake and

Comparison to Tertiary Classification

(a)

(b)

- Investigate whether a distinction between prompt and fake muons offers any advantages for separating them from tau muons
- Loss function: categorical crossentropy $L(y, y^{\text{true}}) = -\sum_{i} y_i^{\text{true}} \log(y_i)$
- Train classifier with flat sampling during batch calls approach, but sample 10,000 prompt as well as fake muons from $t\bar{t}$ during training, so all three classes appear with the same rate
- Consider output score associated with tau muons for all three classes

 \rightarrow Tertiary classification does not improve the performance compared to binary classification

29

Performance Metrics

- Metrics monitored during training: loss and ulletaccuracy
- Normalized output score distributions of all ightarrowthree muon classes from the train and test set
- ROC curve: each point represents a 2D-tupel ightarrowof signal efficiency and background inefficiency for a certain output score threshold; area under curve (AUC) quantifies discrimination power

- Approach: Instead of sampling muons multiple time, use not only muons from $t\bar{t}$, but from all generated processes \bullet
- into train (75%) and test (25%) set, where the latter is defined to be also the evaluation set

250,000 prompt and tau muons are sampled flat in p_T up to 160 GeV. Fakes are not considered here. This dataset is then split up

31

Cause of the intermediate drop can likely be adressed to the correlation of p_T with the impact parameter related input features

- "N-1" feature training: remove one feature during training (in order to resolve strange correlations with p_T)
- Left: p_{T} , $|z_0 \sin \theta|$, $|\eta|$, ρ' were excluded referring to the figures from top to bottom
- Right: $|d_0|$, Δz_0 , $\frac{\Delta p_T}{p_T}$, calorimeter isolation were excluded referring to the figures from top to bottom

- "Shuffle Training": shuffle one feature during training (in order to resolve strange correlations with p_T
- Left: p_{T} , $|z_0 \sin \theta|$, $|\eta|$, ρ' were shuffled referring to the figures from top to bottom
- Right: $|d_0|$, Δz_0 , $\frac{\Delta p_T}{p_T}$, calorimeter isolation were shuffled referring to the figures from top to bottom

• Bin-wise training: trained a DNN for each section (10,40), (40,90), (90,160) GeV individually

(b) Bin-wise training

35

