

# Industrial PhD

*New applications for machine learning and transfer learning algorithms in the manufacturing industry and in particle physics.*

Lorena Dieste





- ✧ Introduction
- ✧ TMVA vs. SKlearn
- ✧ The LHCb Analysis
- ✧ Future work

# Introduction

---





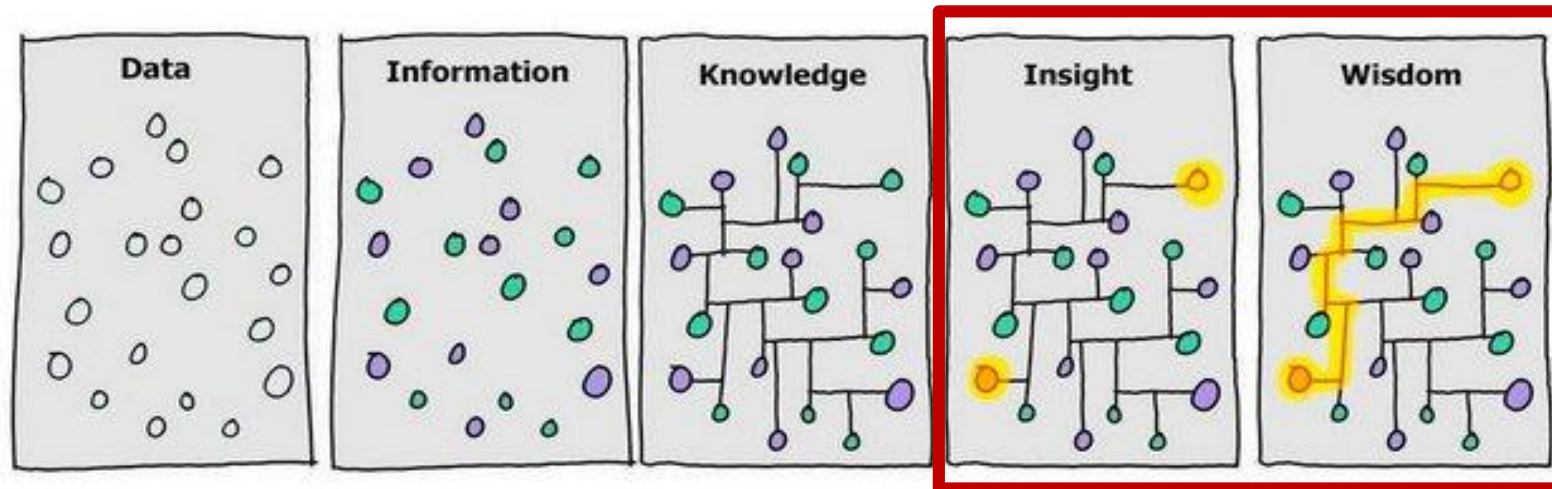
An **industrial PhD student** participates in an industrial research or experimental development project directly related to his thesis, developed in a company and in a university.

- Which company? → TripleAlpha
- **TripleAlpha** is a technological firm specialized in data analysis applied to Production and Operations.
- Which center? → **IGFAE**, LHCb experiment.
- The Galician Institute of High Energy Physics is a joint research center of University of Santiago de Compostela and Xunta de Galicia.



## 1. TripleAlpha

“We help companies to be more productive and profitable thanks to the optimization of production and logistic processes”



Merging **Data Science** with **Manufacturing expertise**





The title of the thesis is:

- *New applications for machine learning and transfer learning algorithms in the manufacturing industry and in particle physics.*

The main goal is to find and use machine learning techniques that can be applied to improve the results for different fields, in the case of this study we will focus on **particle physics** and the **manufacturing industry**.

# TMVA vs. Sklearn

---





## 2. TMVA vs Python

The first step is to compare the tools that are used for particle physics and find the one that gets better results:

- TMVA a ROOT (open-source data analysis framework) tool used by high energy physics and others.
- Python Machine Learning Tools as Scikit-Learn



VS.



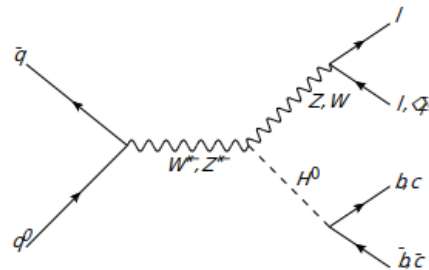
- We will use different types of topologies representing different types of problems of High Energy Physics.





We studied three different types of topologies for this analysis:

- 1)  $K_S^0 \rightarrow \pi^+ \pi^-$  In this exercise our goal is to discriminate signal vs. Background.
- 2)  $D_s \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$  . In this exercise our goal is to discriminate signal vs.. Background.
- 3)  $W+bb$  and  $t\bar{t}$  cross sections. In this exercise our goal is to discriminate them with each other.



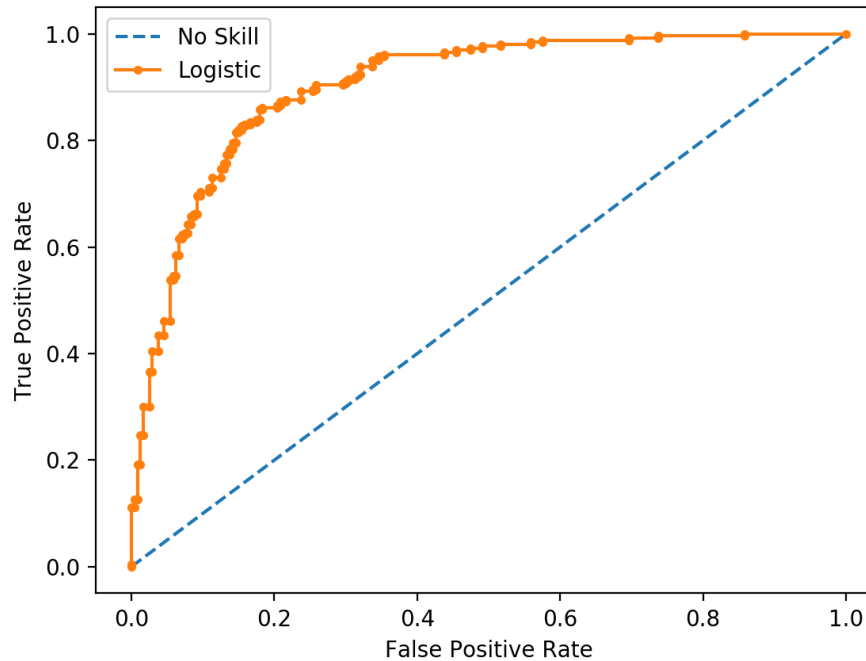
For now we use Monte Carlo LHCb to get the data for the analysis.



## 2. TMVA vs Python

To compare the quality of the results we used two figures of merit:

- The **ROC curve**: a plot that confronts the rate of True Positives vs. False Positives



- The **correlation** between the results and the mass



## 2. TMVA vs Python

To find our model we used **Grid search**. Grid search is used to find the optimal hyperparameters of a model which results in the most 'accurate' predictions.

In our case first we tried different models with different algorithms as:

- BDT with AdaBoost
- BDT with GradientBoosting
- BDT with XGBoost
- RandomForestClassifier
- Neural Network
- Logistic Regression.

```
for c in [2-5, 2-3, ..., 215]:  
    for g in [2-15, 2-13, ..., 25]:  
        for train, test in partition:  
            model = svm_train(train, c, g)  
            score = svm_predict(test, model)  
            cv_list.insert (score)  
        scores_list.insert(mean(cv_list), c, g)  
print max(scores_list)
```

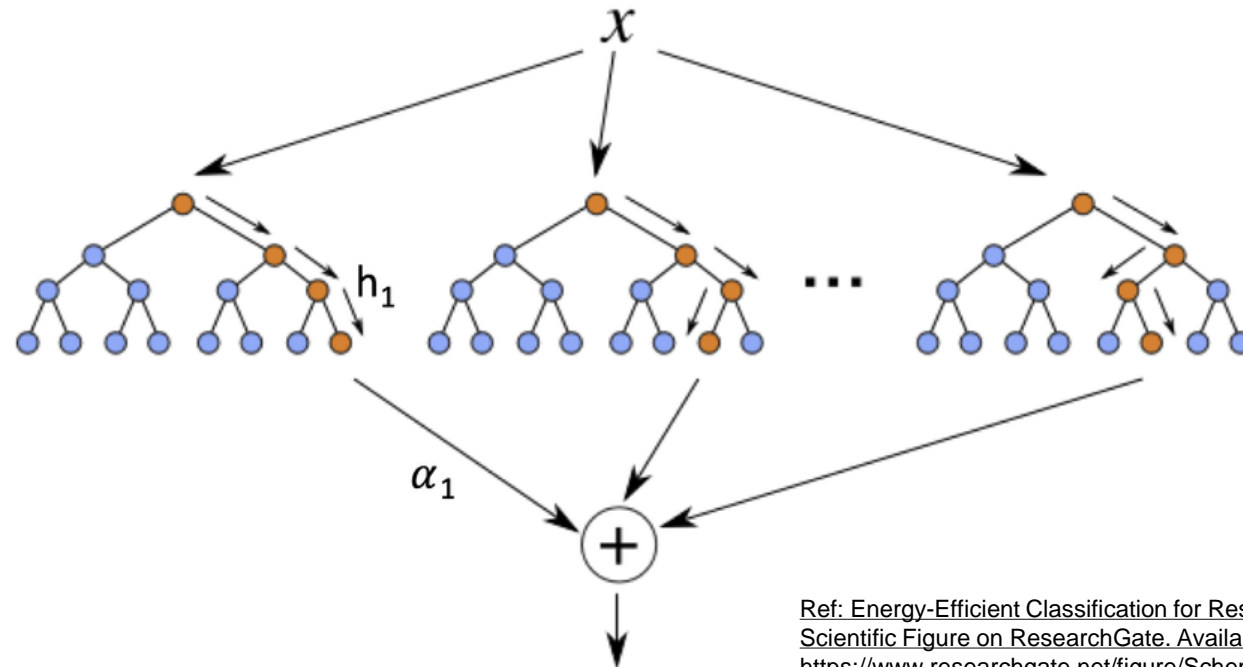
For each one of this we run a Grid Search to find the optimal hyperparameters and after the process we selected the algorithm with the best results.



## 2. TMVA vs Python

The algorithm we decided to use in both platforms was **BDT**, using boosting **AdaBoost**. We also tried Gradient Boosting and XGBoost.

- BDT ( Boosted Decision Trees ) : is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an **ensemble of weak prediction models**, typically decision trees.

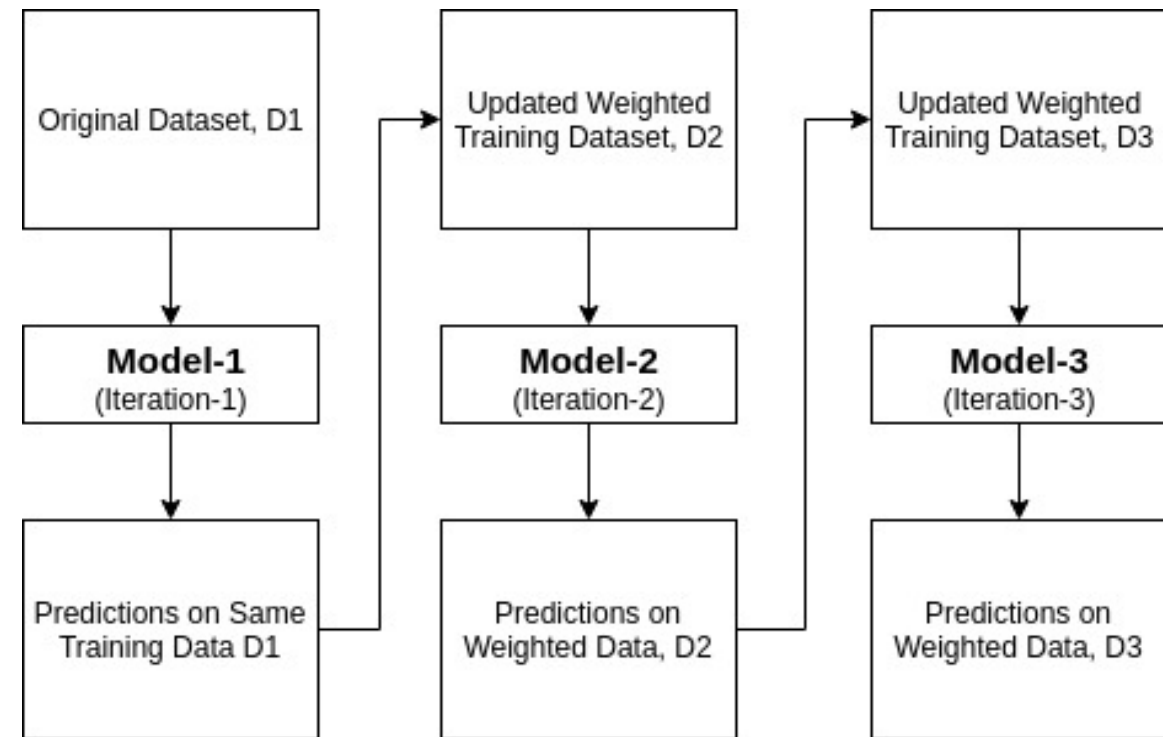




## 2. TMVA vs Python

**Ada-boost** or **Adaptive Boosting** is an iterative ensemble classifier, it builds a strong classifier by combining multiple poorly performing classifiers so that you will get a high accuracy strong classifier.

- It assigns the higher weight to wrong classified observations.
- The more accurate classifier will get high weight.
- This process iterate until the complete training data fits:
  - without any error
  - reached to the specified maximum number of estimators.
- To classify, perform a "vote" across all of the learning algorithms you built.





## 2. TMVA vs Python

---

$K_S^0 \rightarrow \pi^+ \pi^-$  **variables** :

- 'd1\_ip'
- 'd1\_pt'
- 'd2\_ip'
- 'd2\_pt'
- 'doca'
- 'ipchi2'
- 'pt'

Correlation:

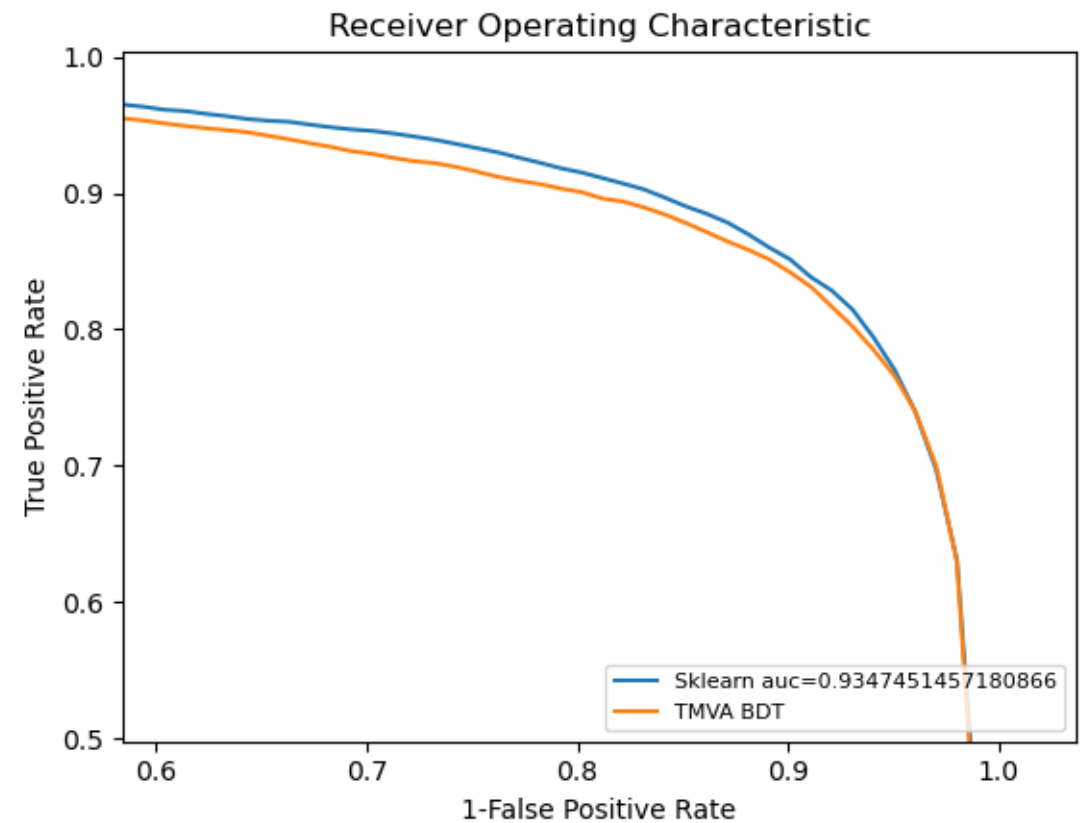
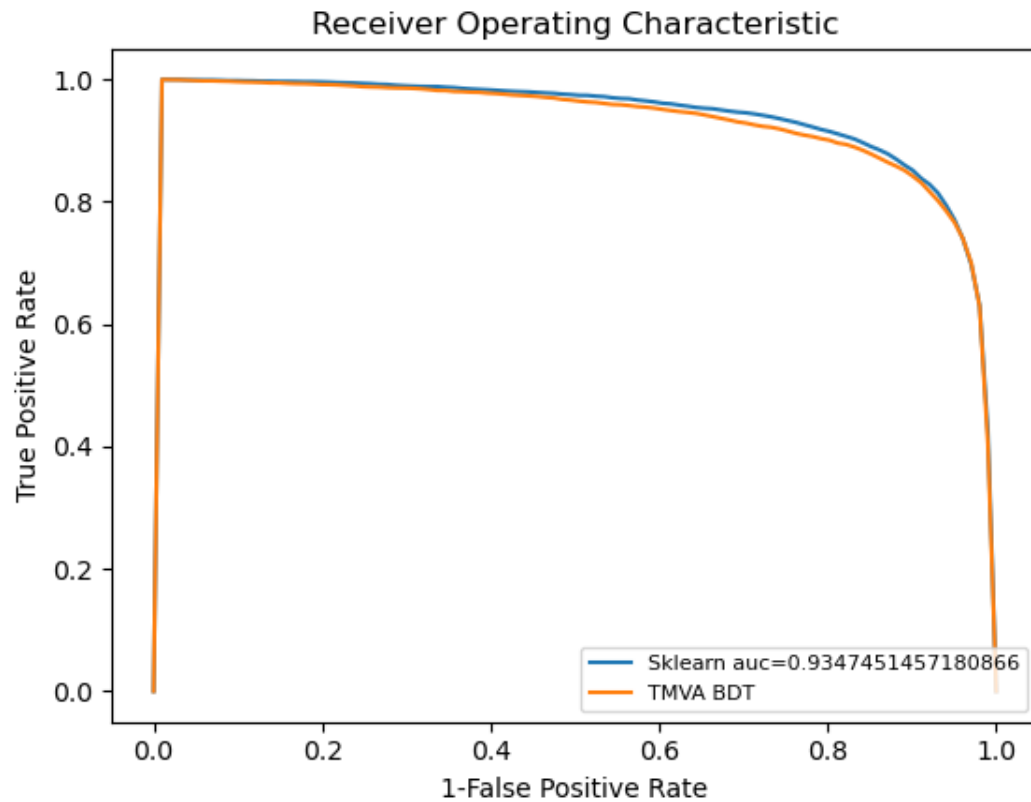
- 'mass'



## 2. TMVA vs Python

KS0  $\rightarrow \pi^+ \pi^-$  results **ROC curve:**

Zoom of the results:

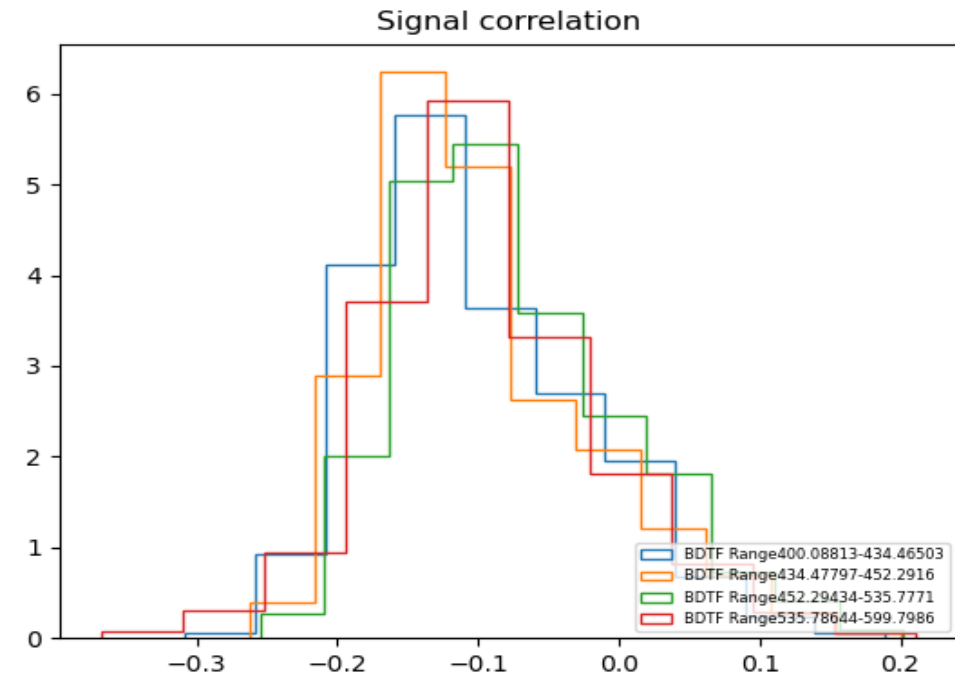
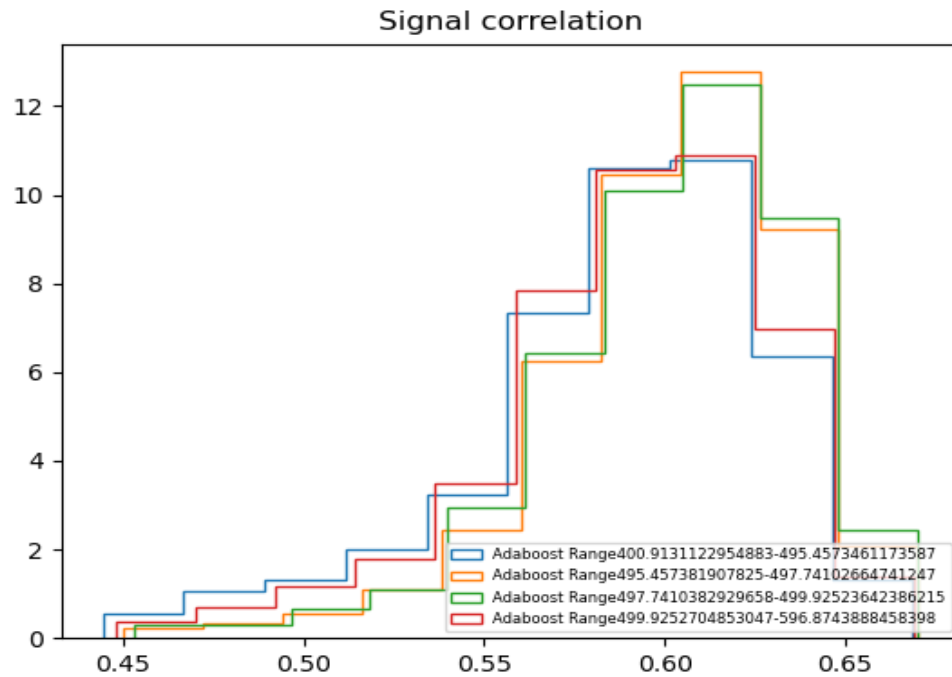


The curve for sklearn is slightly better



## 2. TMVA vs Python

$K_S^0 \rightarrow \pi^+ \pi^-$  results for signal **correlation**:



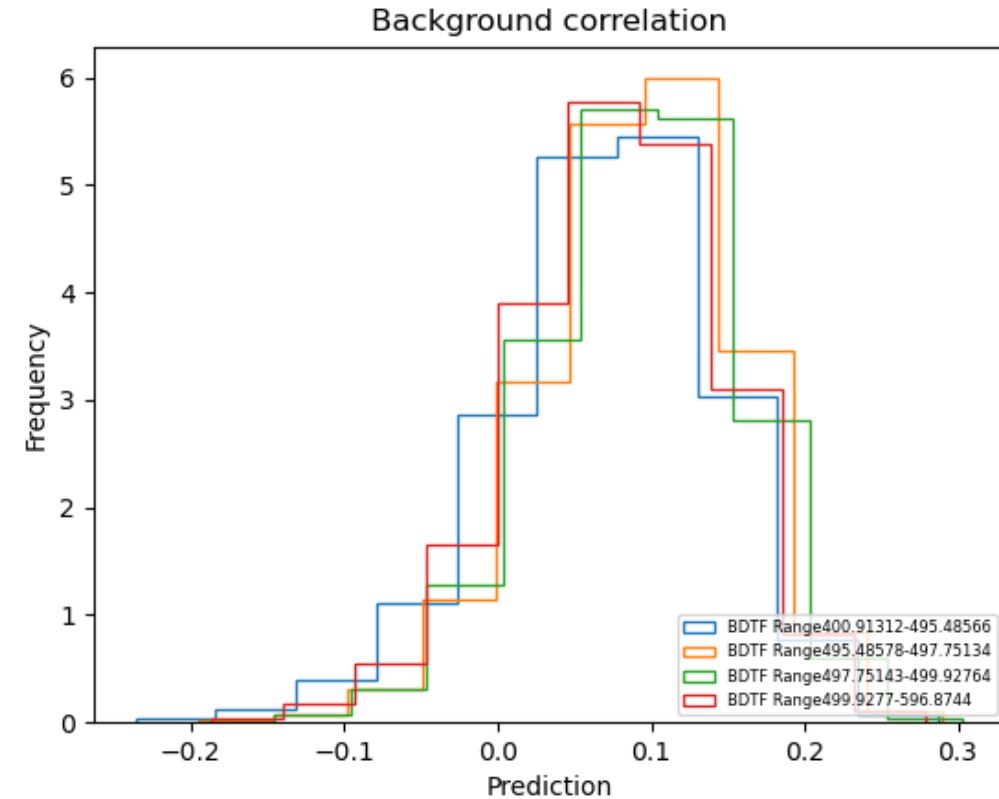
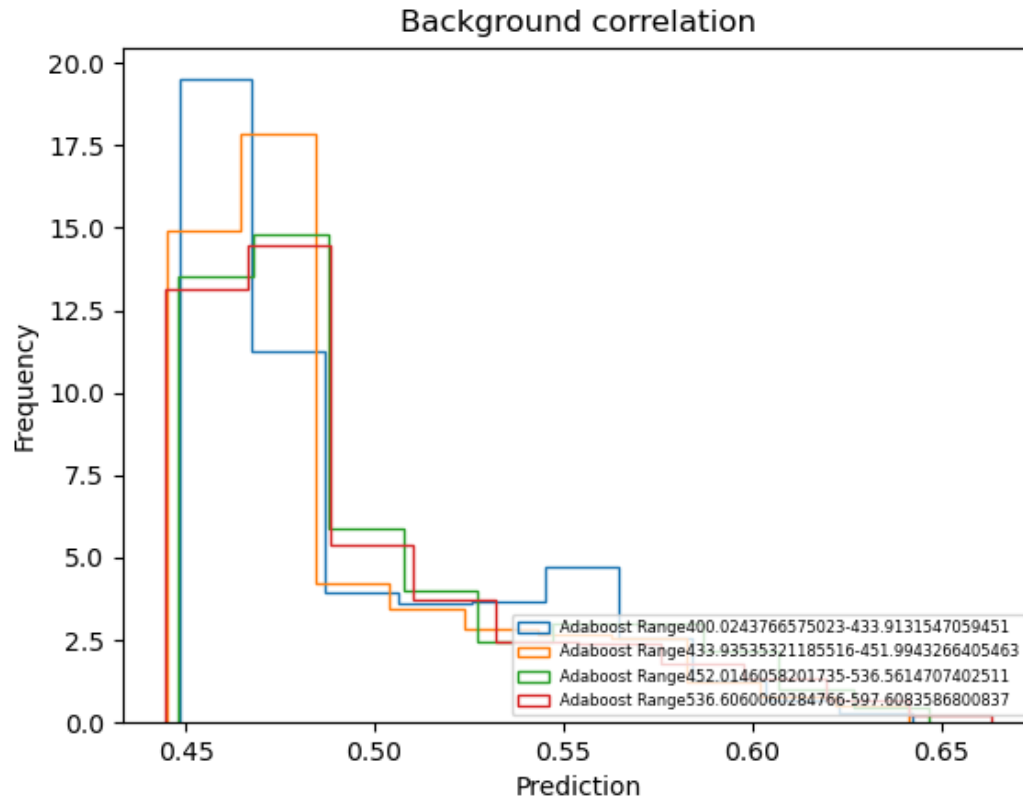
We can observe that there is no correlation.





## 2. TMVA vs Python

$K_S^0 \rightarrow \pi^+ \pi^-$  results for background **correlation**:



We can observe that there is no correlation.



## 2. TMVA vs Python

$Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$  variables :

- 'etap\_IPCHI2\_OWNPV'
- 'etap\_FDCHI2\_OWNPV'
- 'etap\_VCHI2PERDOF'
- 'etap\_PT'
- 'Ds\_IPCHI2\_OWNPV'
- 'Ds\_FDCHI2\_OWNPV',
- 'Ds\_VCHI2PERDOF'
- 'Ds\_PT'
- 'pi\_PT',
- 'pi\_IPCHI2\_OWNPV'
- 'pip\_eta\_PT'
- 'pip\_eta\_IPCHI2\_OWNPV'
- 'pim\_eta\_PT'
- 'pim\_eta\_IPCHI2\_OWNPV'
- 'mup\_PT'
- 'mup\_IPCHI2\_OWNPV'
- 'mum\_PT'
- 'mum\_IPCHI2\_OWNPV']

Correlation:

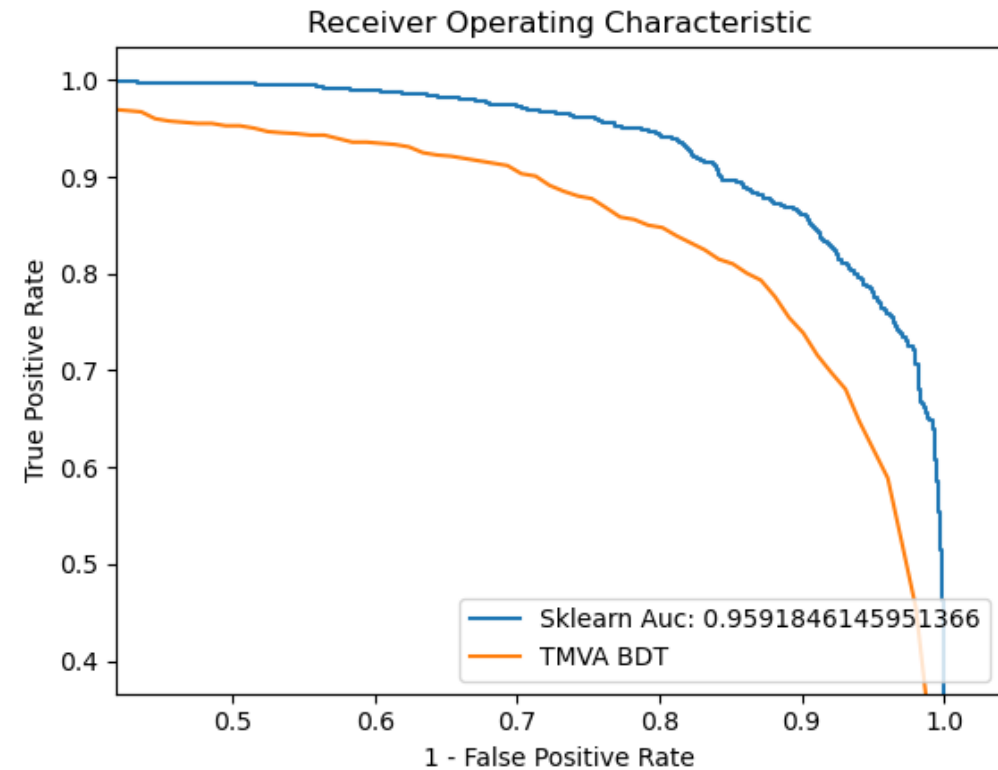
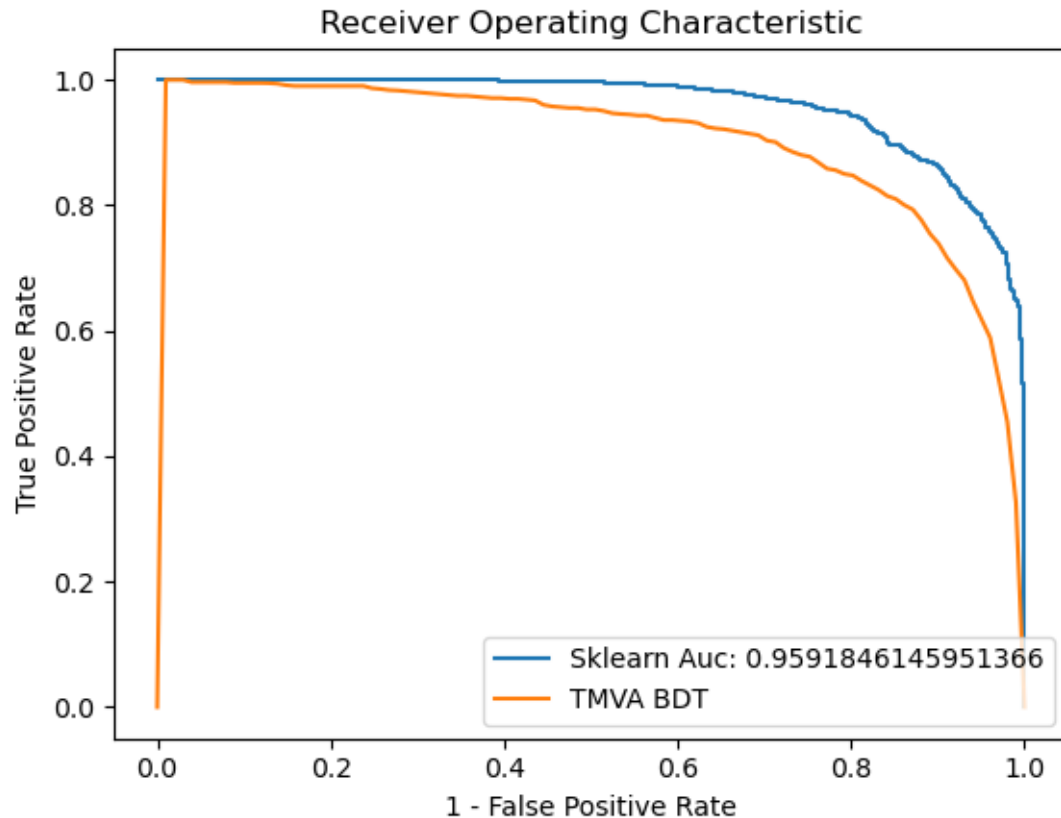
- 'Ds\_M'
- 'etap\_M'



## 2. TMVA vs Python

In  $D_s \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$  we used a K-fold to be able to use all the data for training, the results of the **ROC curve**:

Zoom of the results:

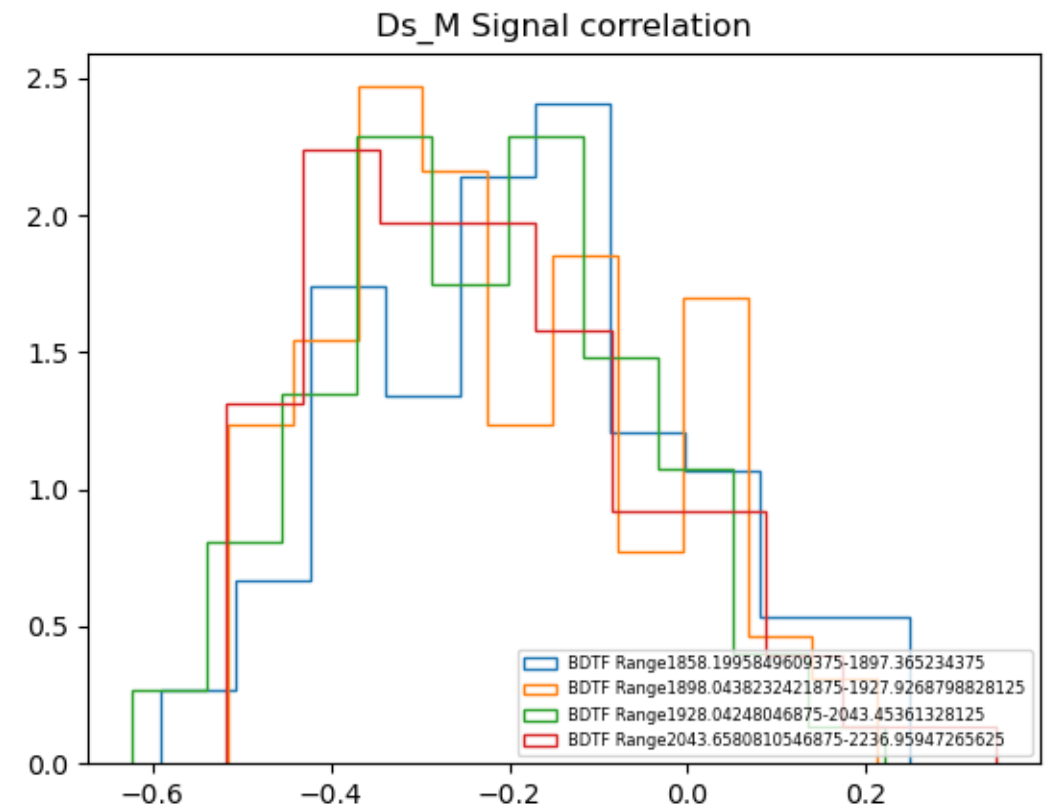
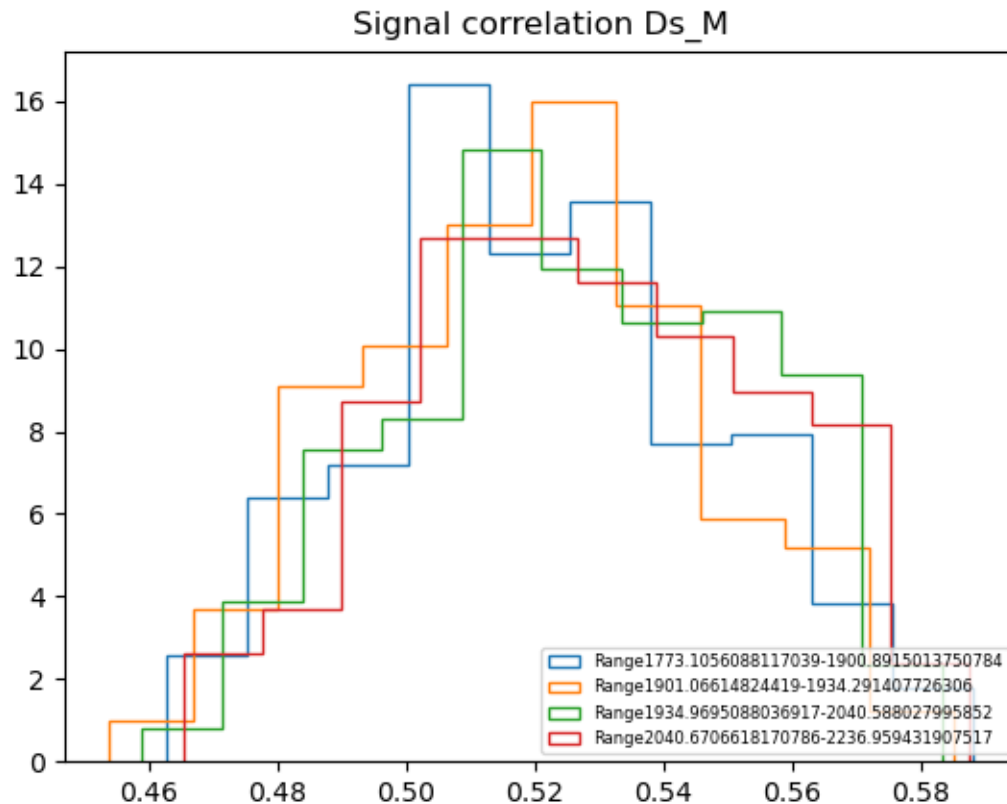


The curves for sklearn have a visible improvement



## 2. TMVA vs Python

Second experiment results for signal **correlation**:

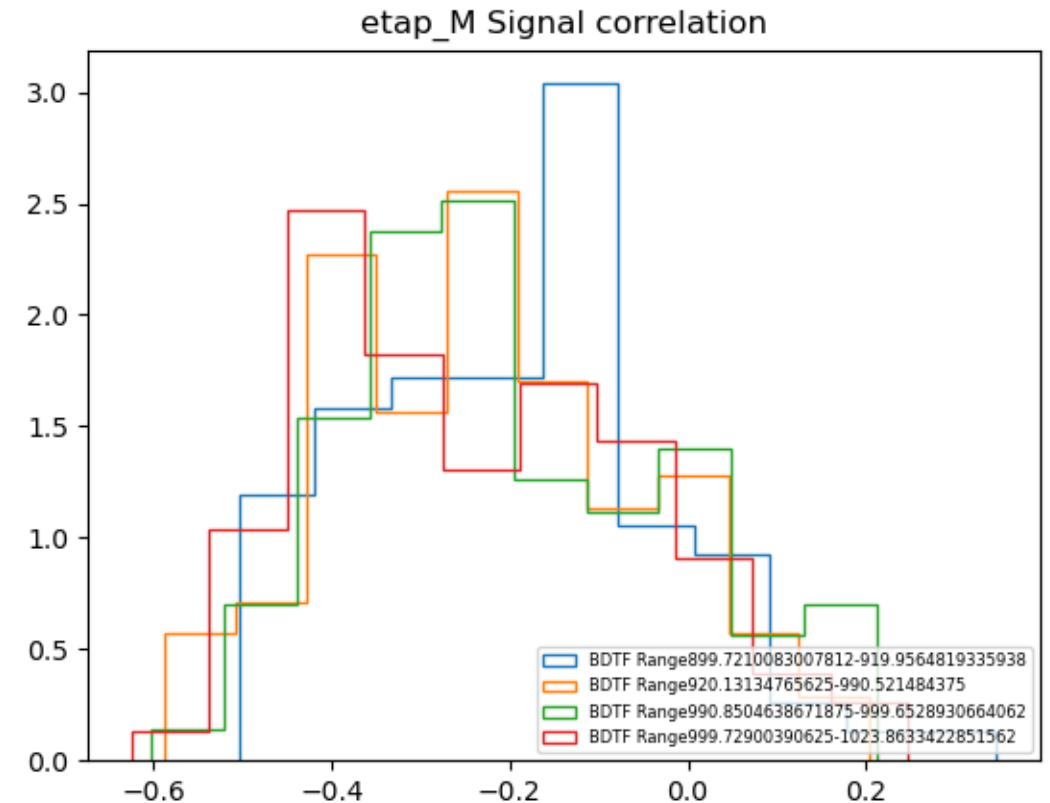
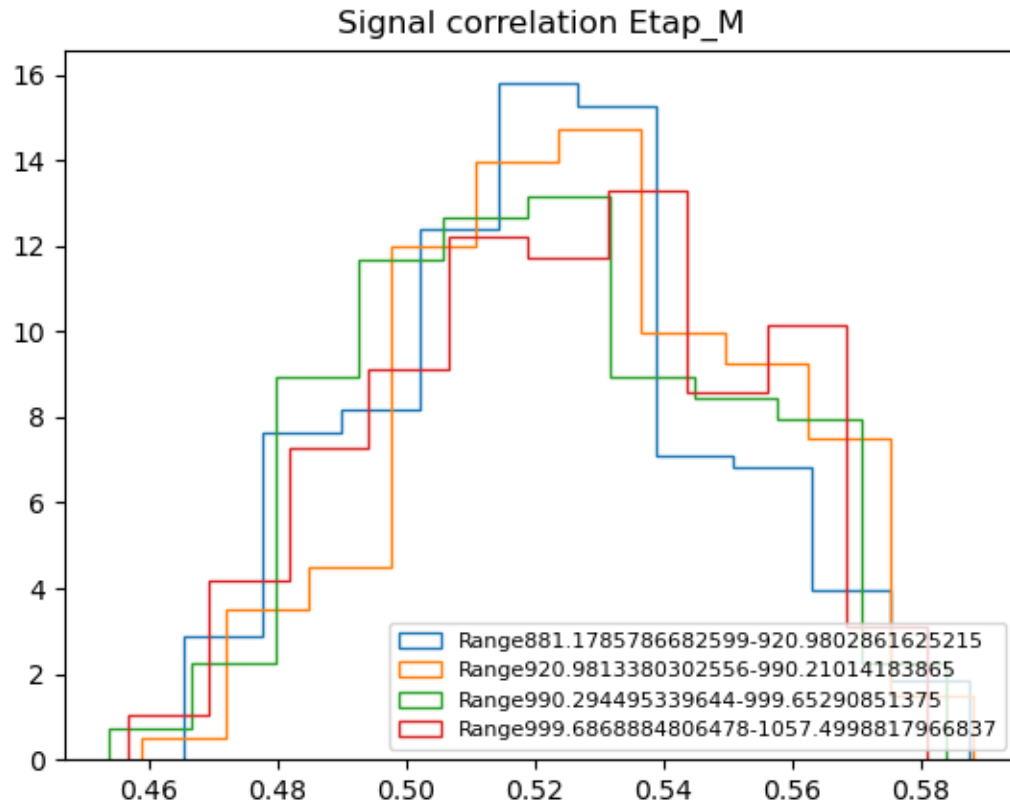


We can observe that there is no correlation.



## 2. TMVA vs Python

In this example the particle decays from one corps to another and we have two masses, Ds\_M and Etap\_M

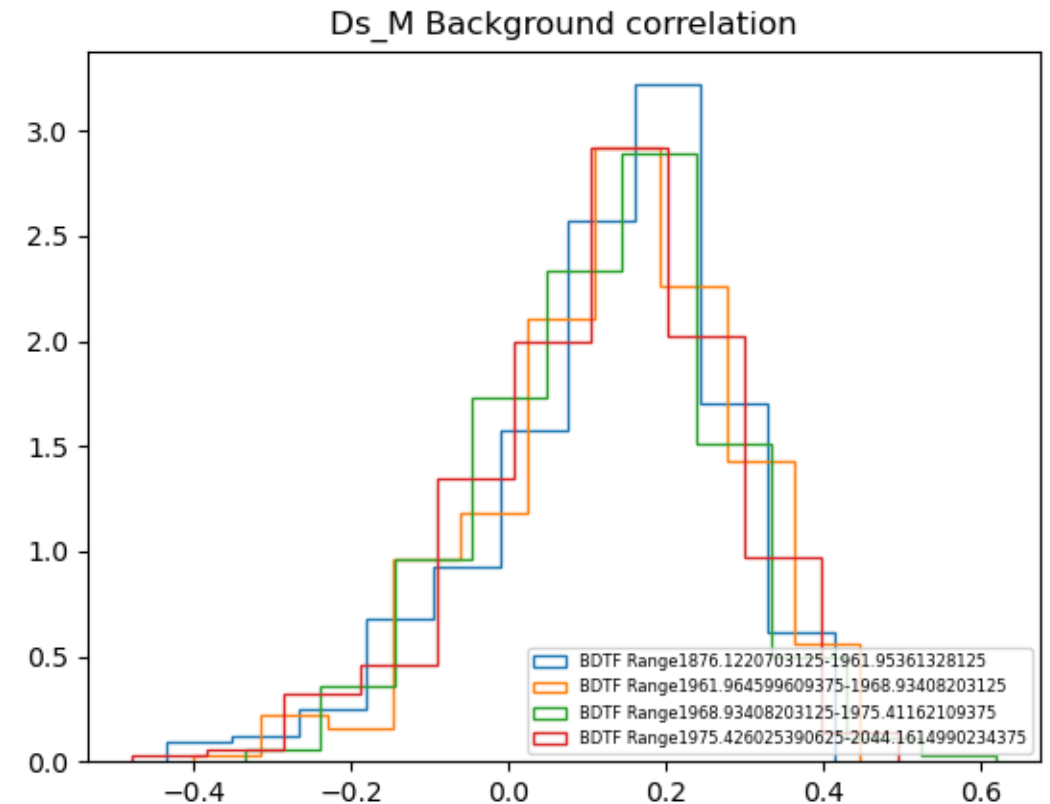
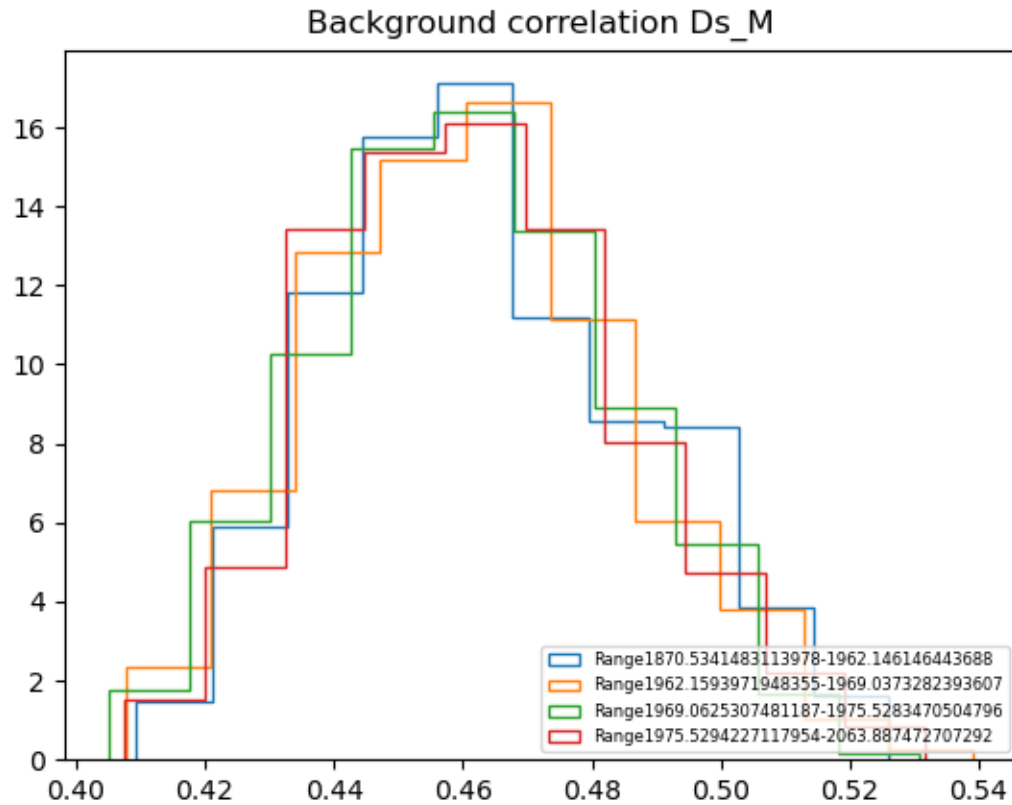


We can observe that there is no correlation.



## 2. TMVA vs Python

Second experiment results for background **correlation**:

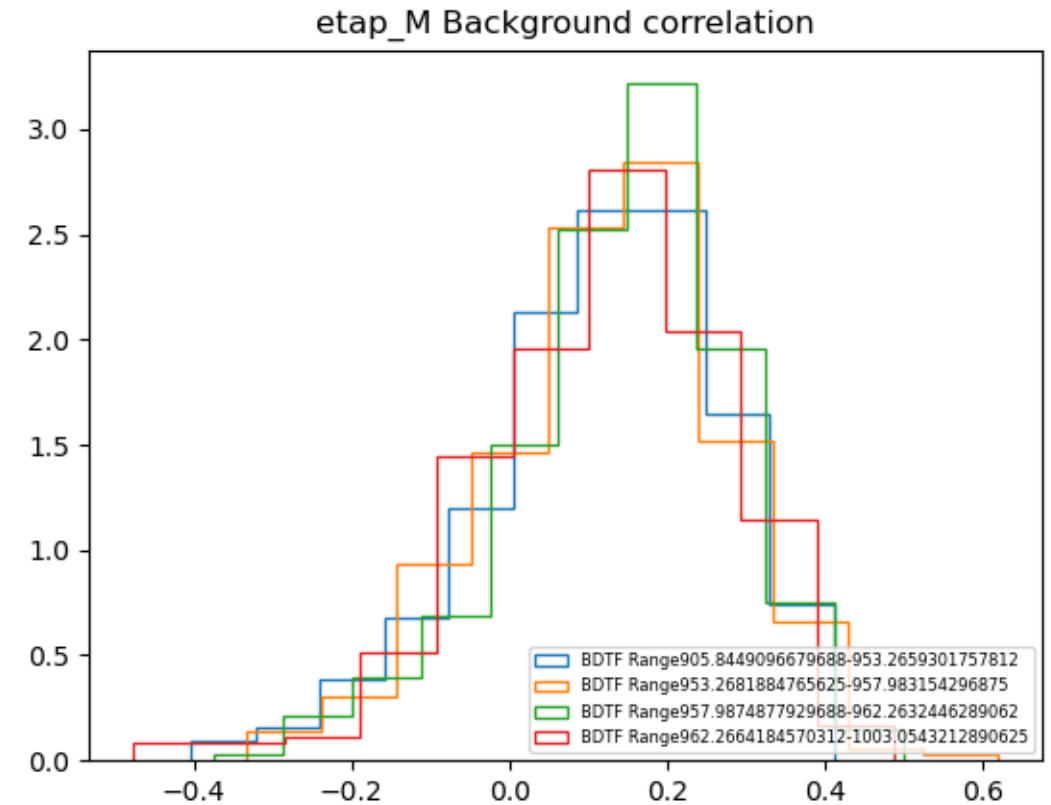
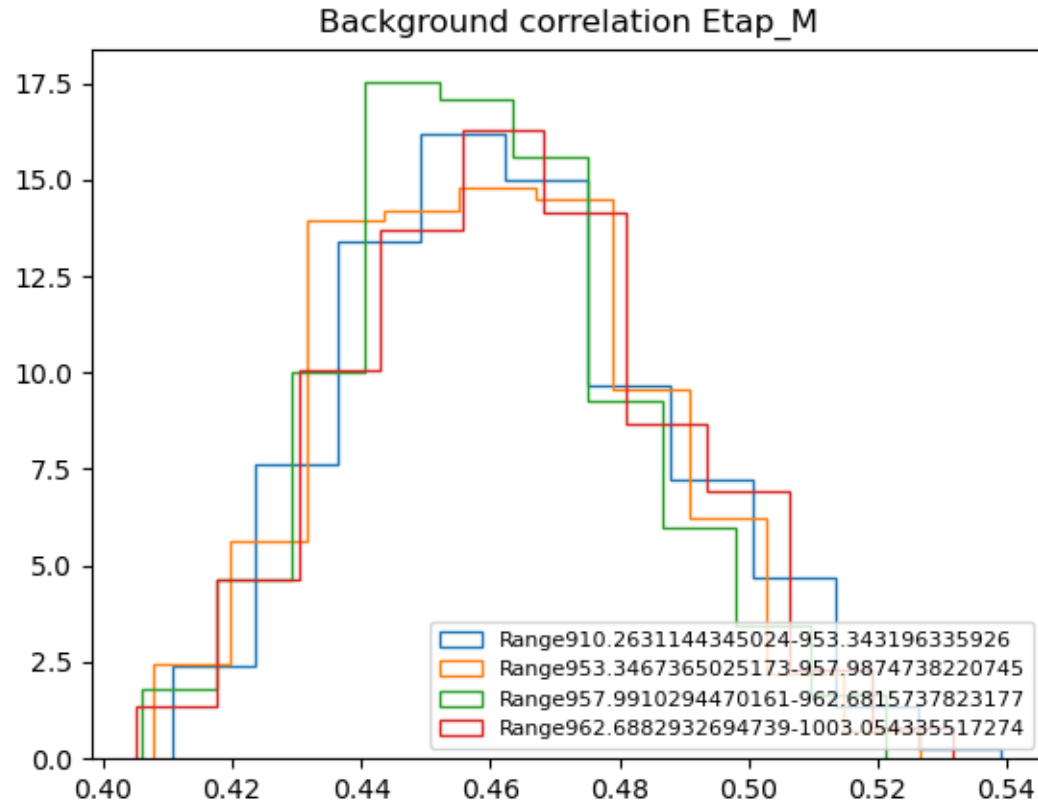


We can observe that there is no correlation.



## 2. TMVA vs Python

Second experiment results for background **correlation**:



We can observe that there is no correlation.



## 2. TMVA vs Python

---

W+bb and tt :

- "Lepton\_ETA"
- "Lepton\_PT"
- "DiJet\_PT"
- "HP\_DR\_JET0\_LEPTON"
- "HP\_DR\_JET1\_LEPTON"
- "Jet0\_JetWidth"
- "Jet1\_JetWidth"
- "HP\_CosThetaBoost"
- "HP\_DR\_JET0\_JET1"
- "ULE\_CELLjet0\_PT"
- "PF\_ETA\_VAR\_ERatio\_12\_3456"

Correlation:

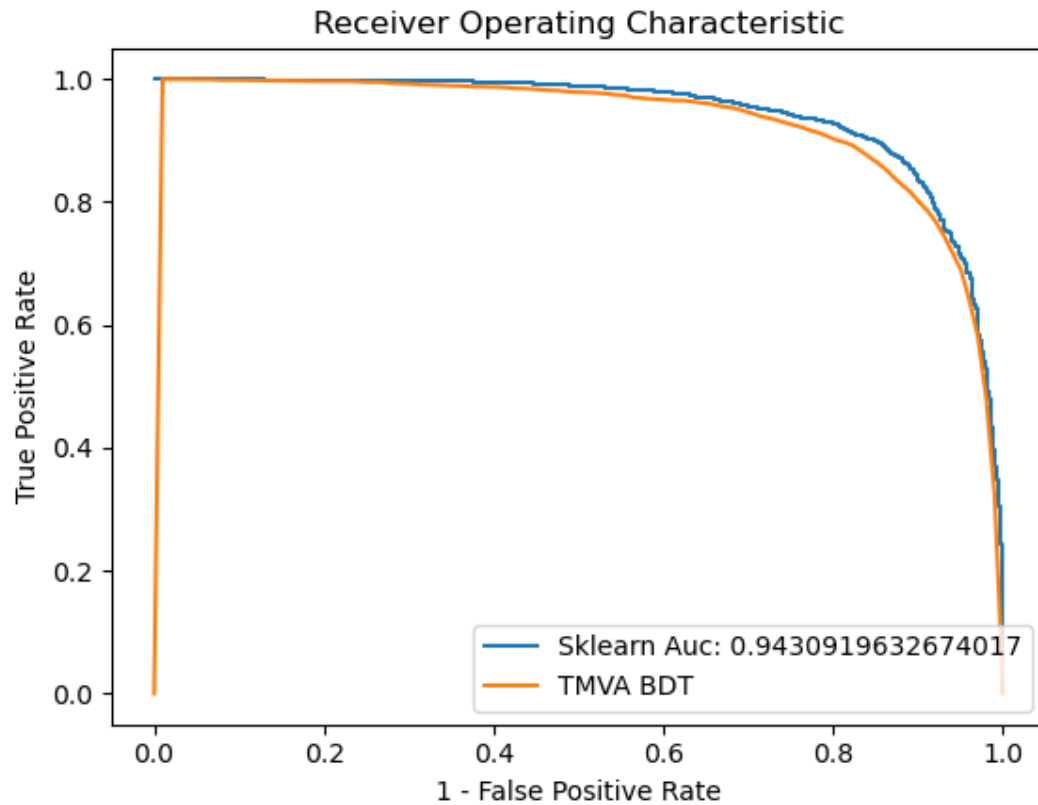
- "DiJet\_M"



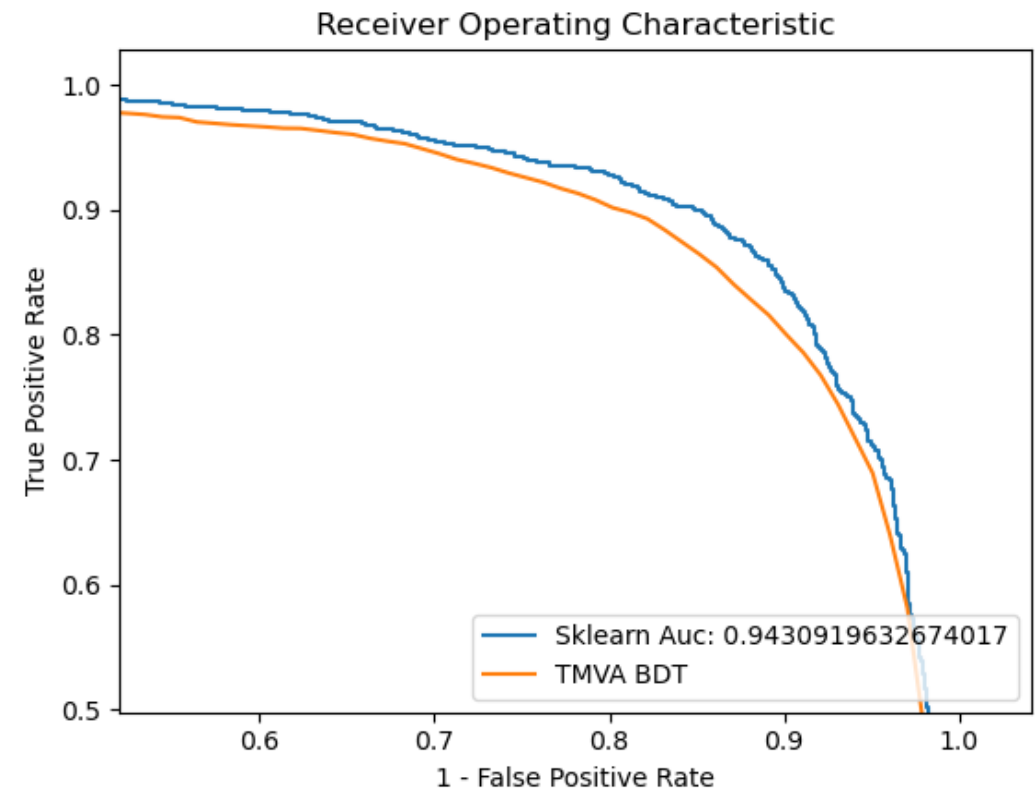


## 2. TMVA vs Python

W+bb and tt results **ROC curve:**



Zoom of the results:

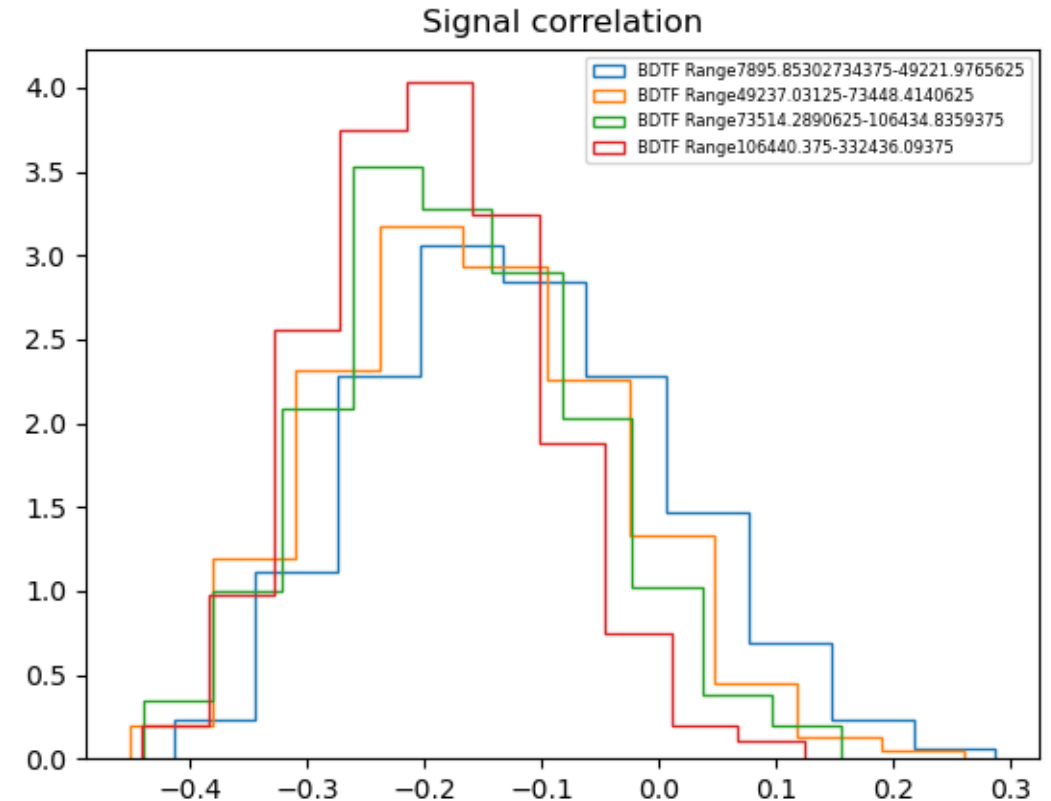
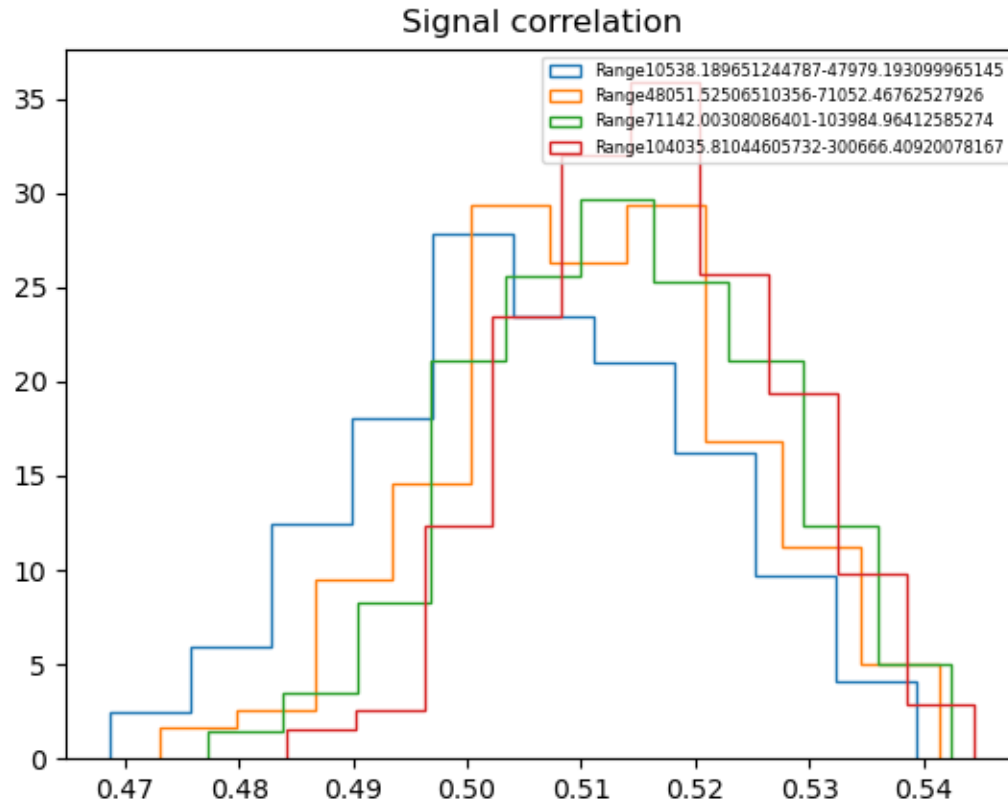


The curve for sklearn is better than the TMVA's



## 2. TMVA vs Python

Third experiment results for signal **correlation**:

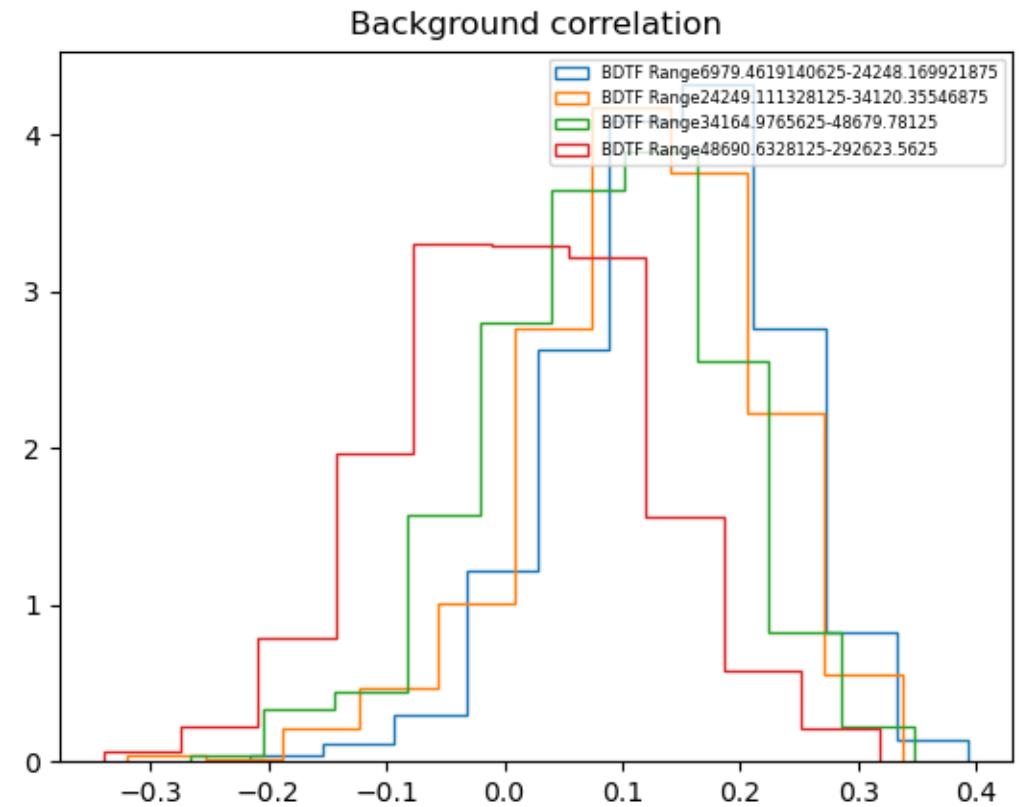
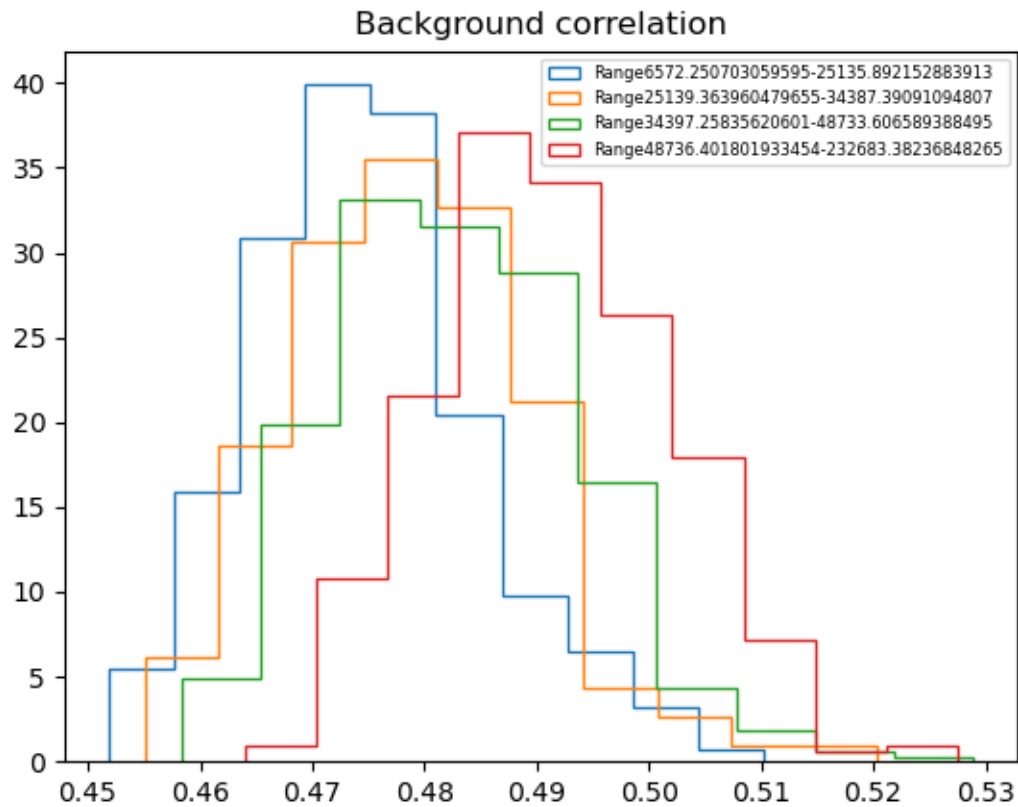


We can observe that there is no correlation in both results.



## 2. TMVA vs Python

Third experiment results for background **correlation**:



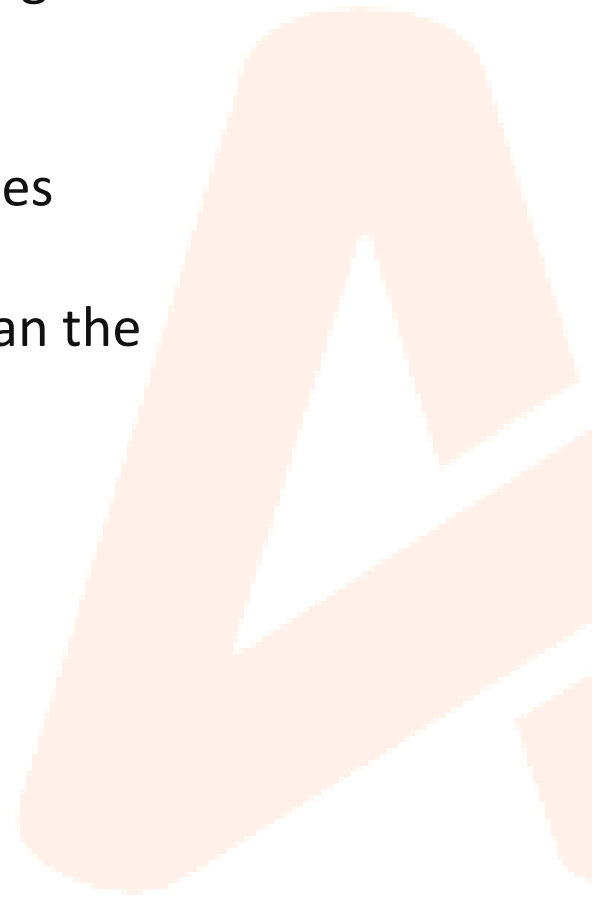
We can observe that there is correlation in both results.



## 2. Conclusions

---

- After three different results we concluded that the best results were available using Scikit-Learn instead of TMVA.
- The ROC Curves in all cases were better for the scikit-learn results using the same parameters
- The correlation with the mass behaved similarly enough in both cases
- In conclusion the results obtained using Scikit-learn are preferred than the TMVA ones for the same problem and same specifications.





## 2. Future Work

---

- 1) We need to keep proving with different data that this conclusion is correct and not just a singular occasion in this specific data.
- 2) We will use Pythia to generate data for the next analysis in order to avoid the use of LHCb data that can't be used for publications outside the CERN.
- 3) We need to prove why the results under the same specifications work better for our module in scikit-learn Python than our module in TMVA Python.



# LHCb Analysis

---





Another line of work of the thesis is to help in different analysis in the LHCb collaboration.

We got three different groups of data:

- $Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$
- $Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$
- $Ds \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$

In this case we want to measure the branching fractions of the decays that are theoretically predicted, this experimental measurement is useful for understanding Chiral Perturbation Theory. [1]

1. GAN, L., KUBIS, B., PASSEMAR, E. Y TULIN, S.

Precision tests of fundamental physics with  $\eta$  and  $\eta'$  mesons



Our goal was to create a Python module, using the same Machine Learning algorithm as before (BDT), to get the results of the classification for the three different types of data:

- $Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$
  - $Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$
  - $Ds \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$
- We used as data the LHCb RUN2 dataset including 2016, 2017 and 2018. As the data we had was not extensive we decide to use a K-fold of the different years as training for the others. We compared the variables of the different years to prove that the data was similar enough to train the models.





$Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$  variables :

- 'etap\_IPCHI2\_OWNPV'
- 'etap\_FDCHI2\_OWNPV'
- 'etap\_VCHI2PERDOF'
- 'etap\_PT'
- 'Ds\_IPCHI2\_OWNPV'
- 'Ds\_FDCHI2\_OWNPV',
- 'Ds\_VCHI2PERDOF'
- 'Ds\_PT'
- 'pi\_PT',
- 'pi\_IPCHI2\_OWNPV'
- 'pip\_eta\_PT'
- 'pip\_eta\_IPCHI2\_OWNPV'
- 'pim\_eta\_PT'
- 'pim\_eta\_IPCHI2\_OWNPV'
- 'mup\_PT'
- 'mup\_IPCHI2\_OWNPV'
- 'mum\_PT'
- 'mum\_IPCHI2\_OWNPV']

Correlation:

- 'Ds\_M'
- 'etap\_M'

$Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi'$  :

- 'etap\_IPCHI2\_OWNPV'
- 'etap\_FDCHI2\_OWNPV'
- 'etap\_VCHI2PERDOF'
- 'etap\_PT'
- 'Ds\_IPCHI2\_OWNPV'
- 'Ds\_FDCHI2\_OWNPV'
- 'Ds\_VCHI2PERDOF'
- 'Ds\_PT'
- 'pi\_PT'
- 'pi\_IPCHI2\_OWNPV'

Correlation:

- 'Ds\_M'
- 'etap\_M'



$Ds \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$ :

- 'eta\_IPCHI2\_OWNPV'
- 'eta\_FDCHI2\_OWNPV'
- 'eta\_VCHI2PERDOF'
- 'eta\_PT','Ds\_IPCHI2\_OWNPV'
- 'Ds\_FDCHI2\_OWNPV'
- 'Ds\_VCHI2PERDOF'
- 'Ds\_PT'
- 'pi\_PT'
- 'pi\_IPCHI2\_OWNPV'

Correlation:

- 'Ds\_M'
- 'eta\_M'

Isolation for  $Ds \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$

And  $Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$ :

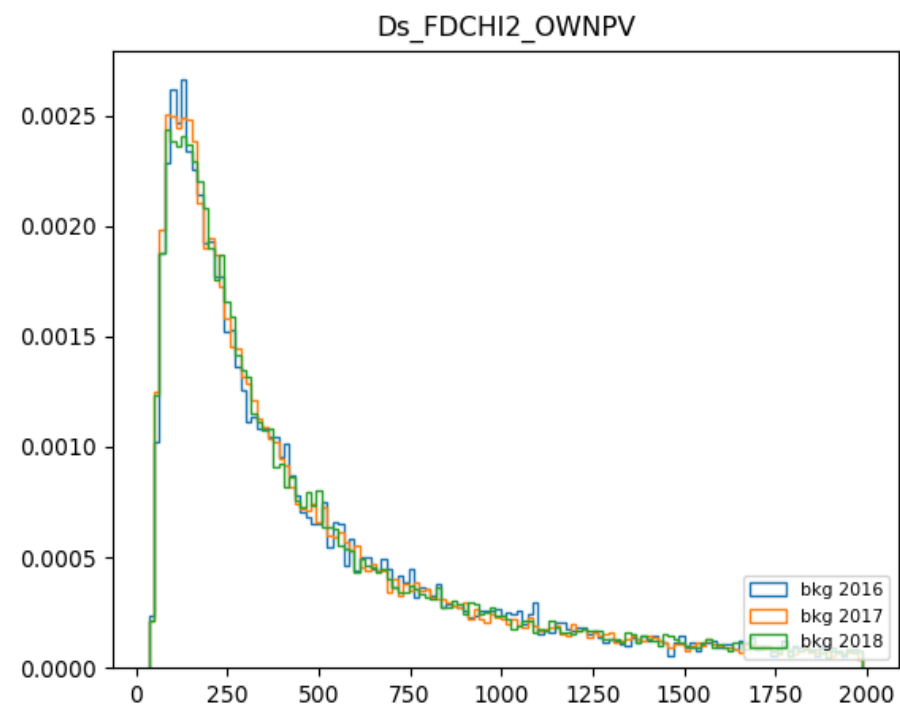
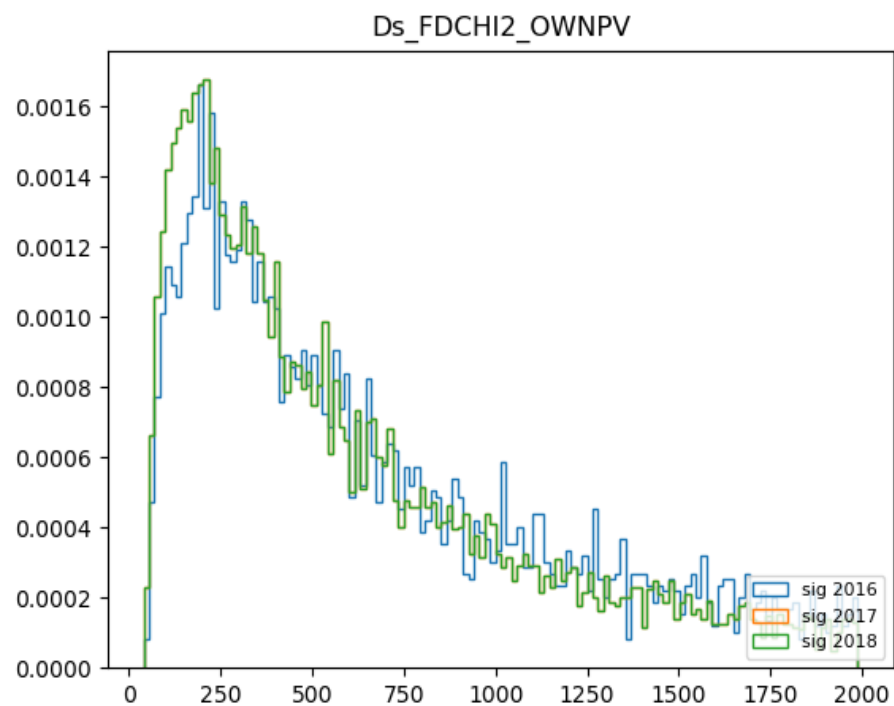
- 'Ds\_VertexIsolation\_VTXISODCHI2ONETRACK'
- 'Ds\_VertexIsolation\_VTXISODCHI2TWOTRACK'
- 'pi\_TrackIsolationBDT\_TRKISOBDDTTHIRDVALUE'

Isolation for  $Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$

- 'Ds\_VertexIsolation\_VTXISODCHI2ONETRACK'
- 'Ds\_VertexIsolation\_VTXISODCHI2TWOTRACK'

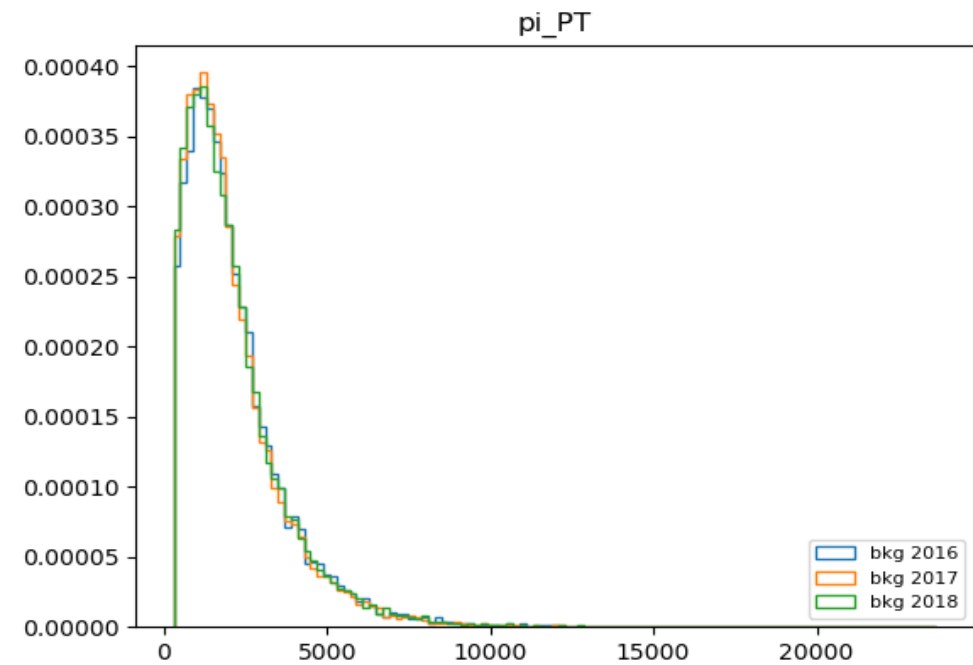
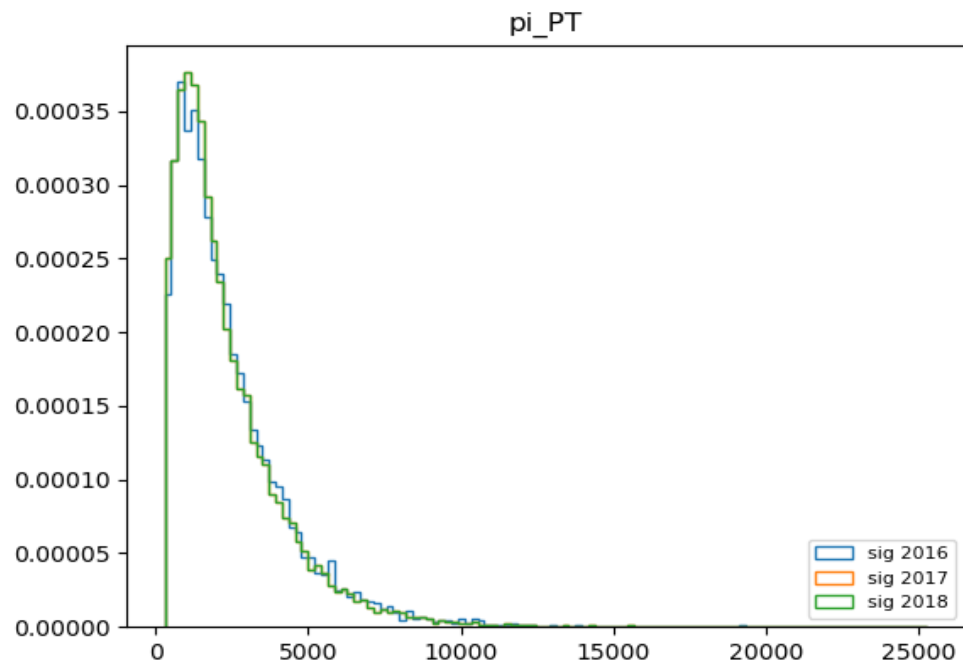


$$D_s \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$$



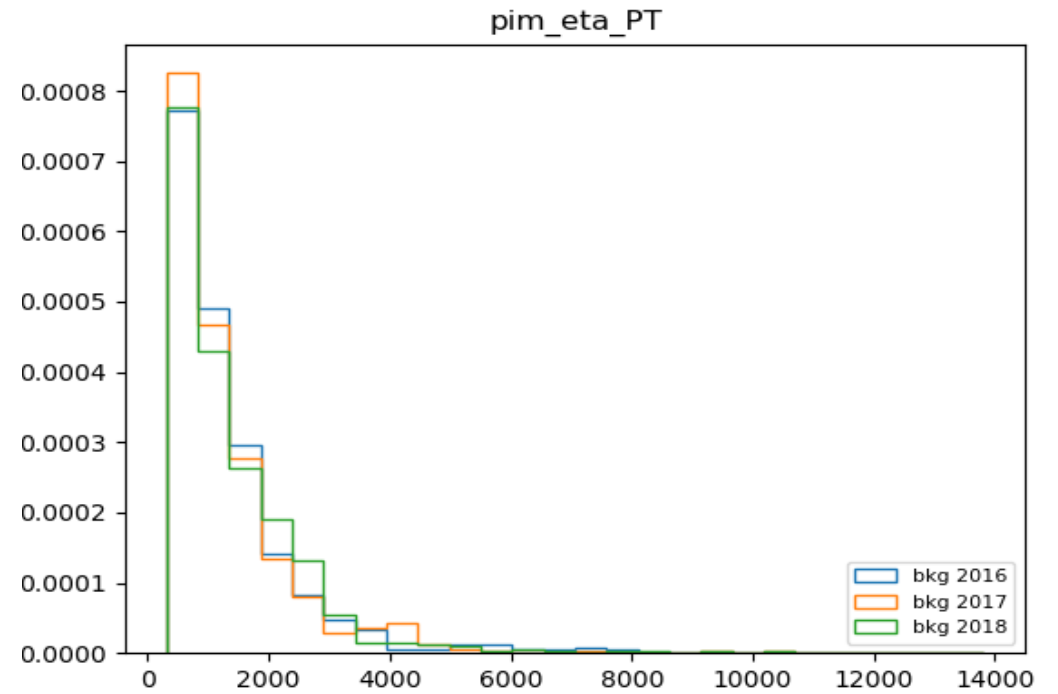
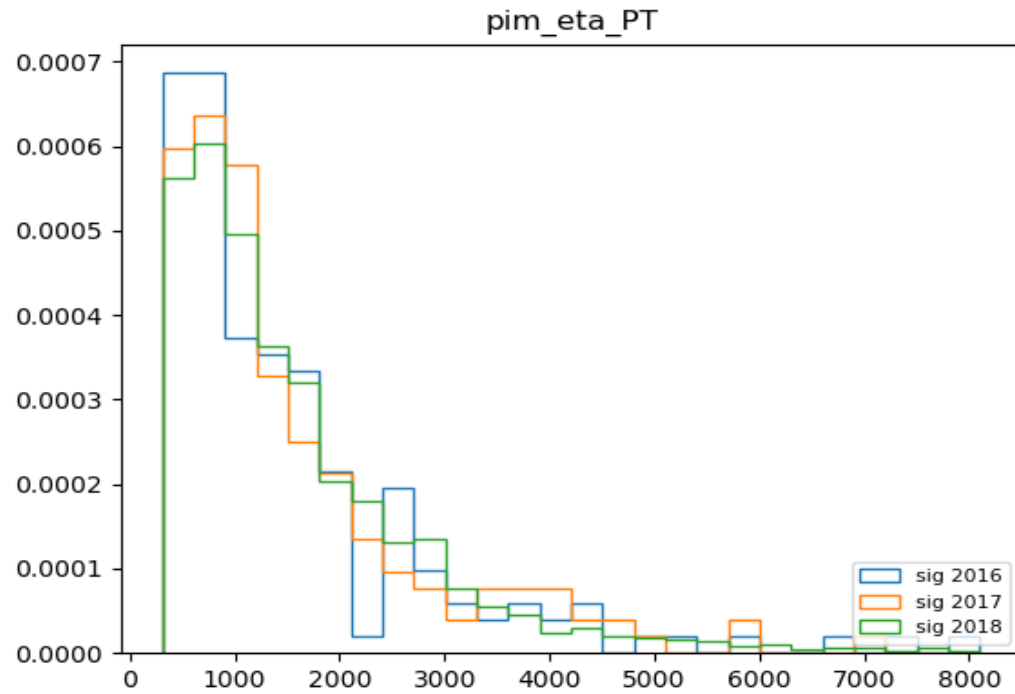


- $Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$





$$D_s \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$$





To measure the quality of the results we used same options as before:

1)The **ROC curve**: a plot that confronts the rate of True Positives vs. False Positives

1)The **correlation** between the results and the mass

We got all the results for each analysis and each year (2016, 2017, 2018), in total 9 different classifiers.

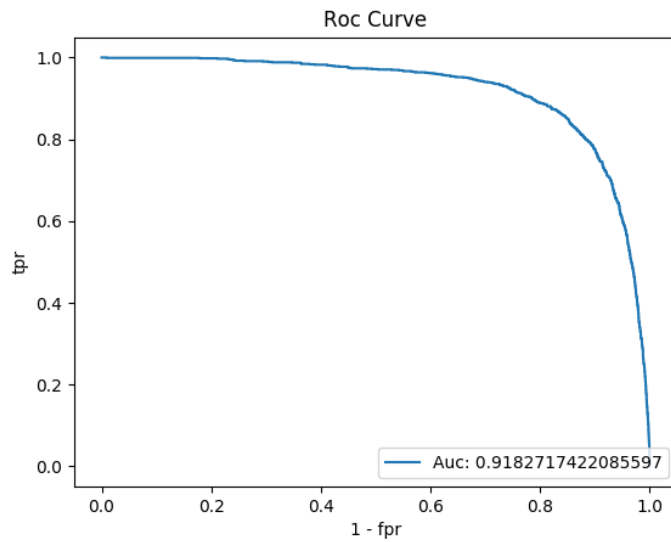
```
1. Ds->eta->(mumu)+pi
2. Ds->etap->(mumu)+pi
3. Ds->etap(->mumupipi)+pi
4. Exit/Quit
Option: 
```



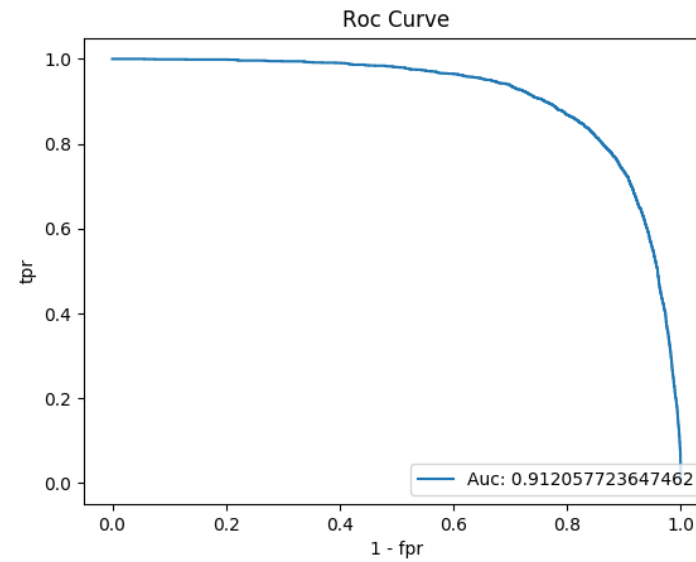
### 3. Results

$$D_s \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$$

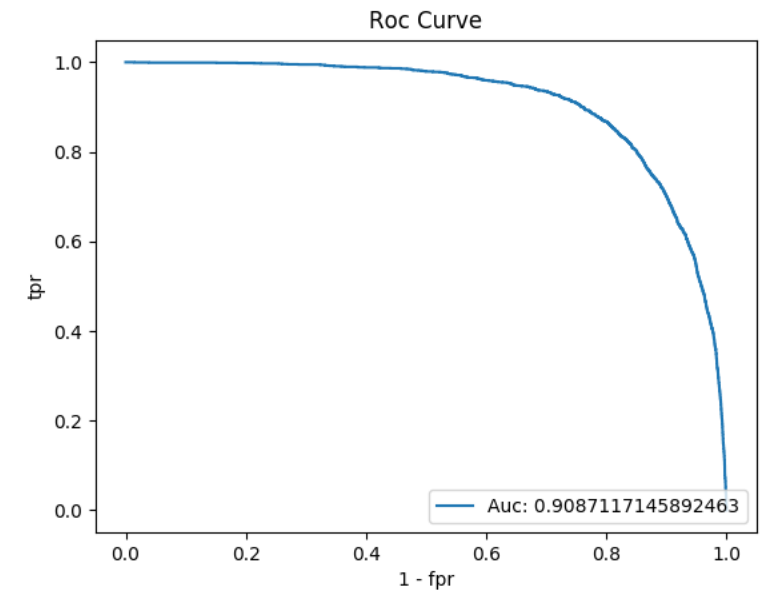
2016



2017



2018

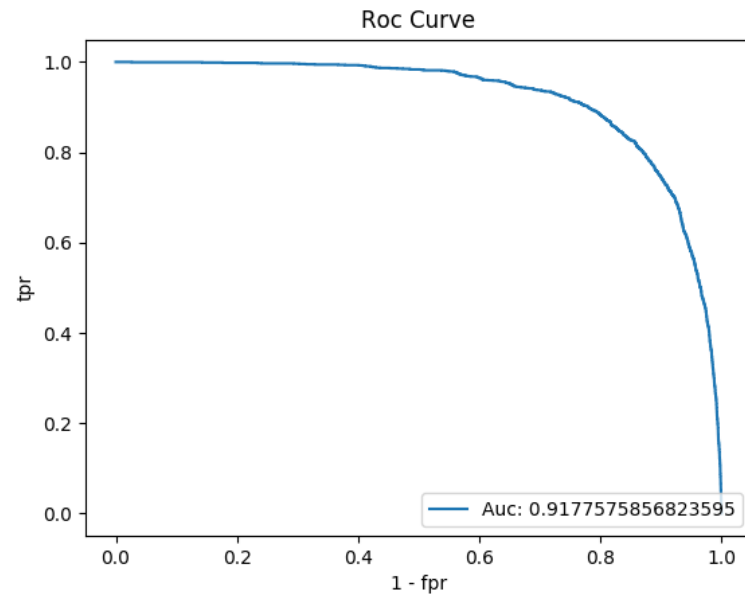




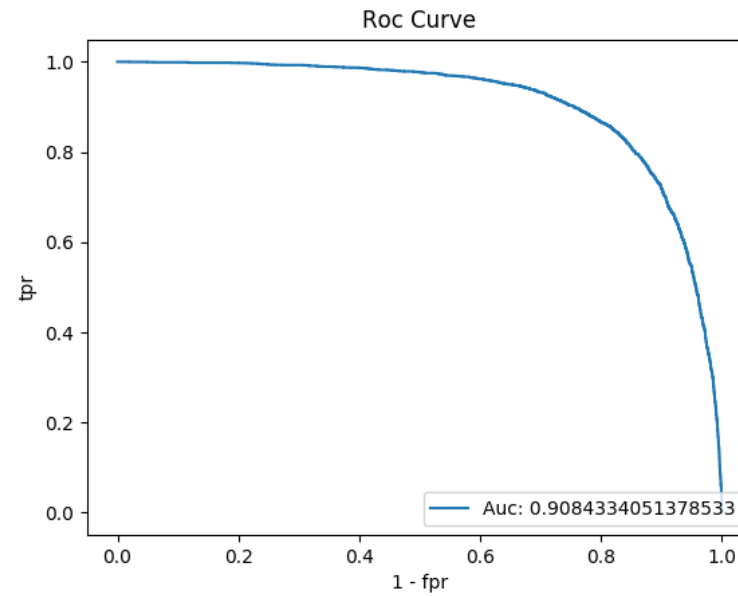
### 3. Results

$$D_s \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$$

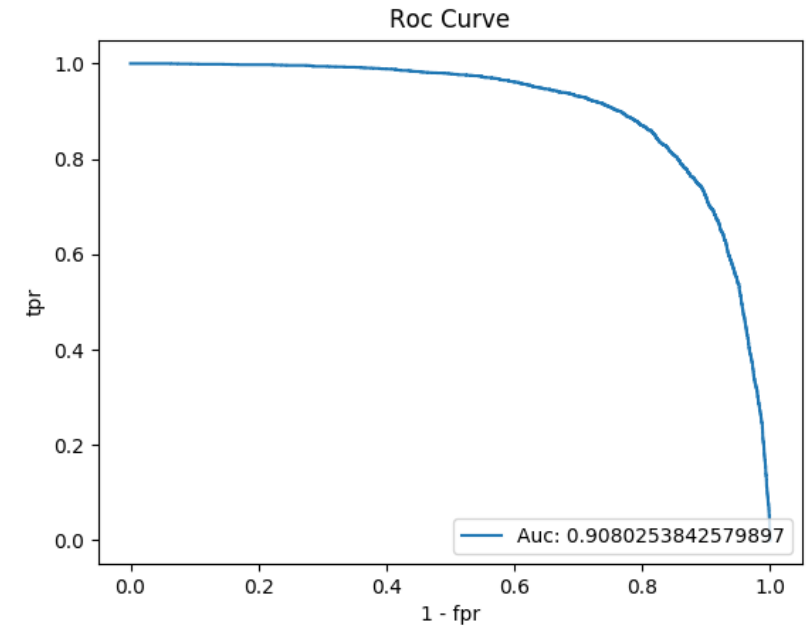
2016



2017



2018



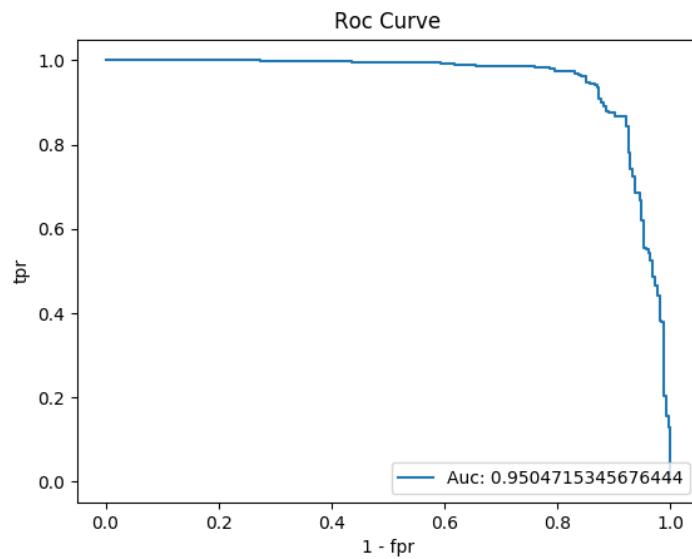




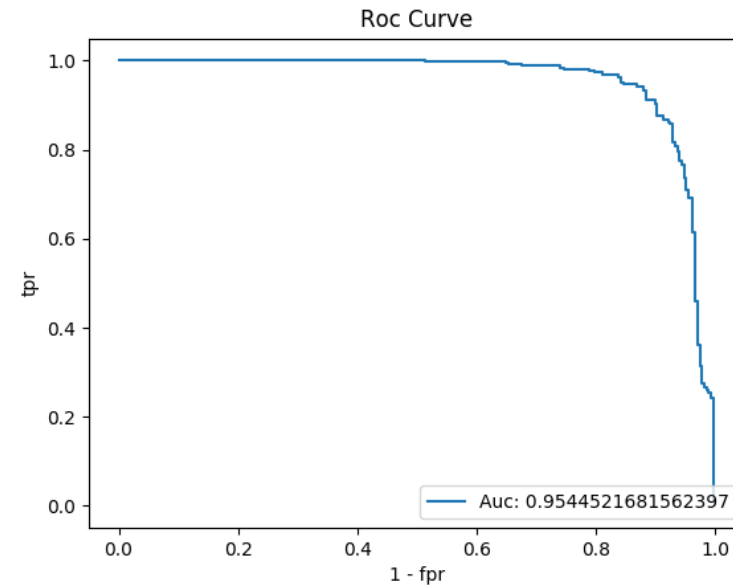
### 3. Results

$$D_s \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$$

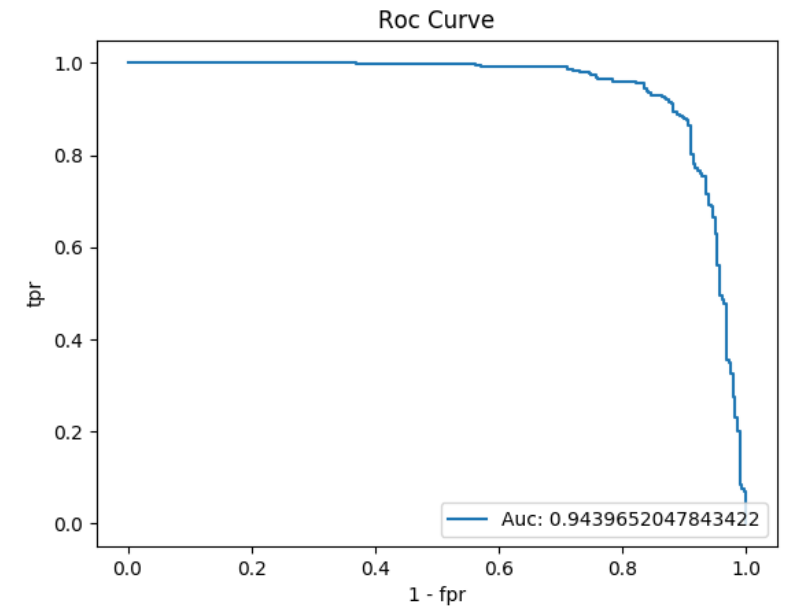
2016



2017



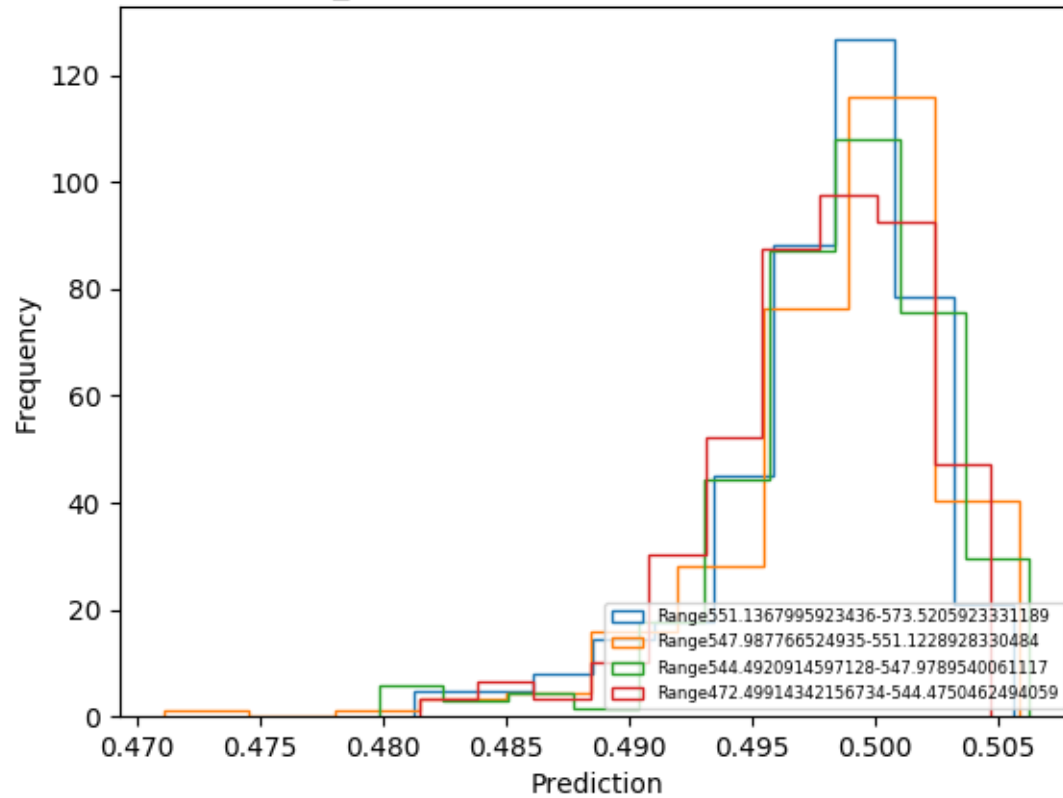
2018



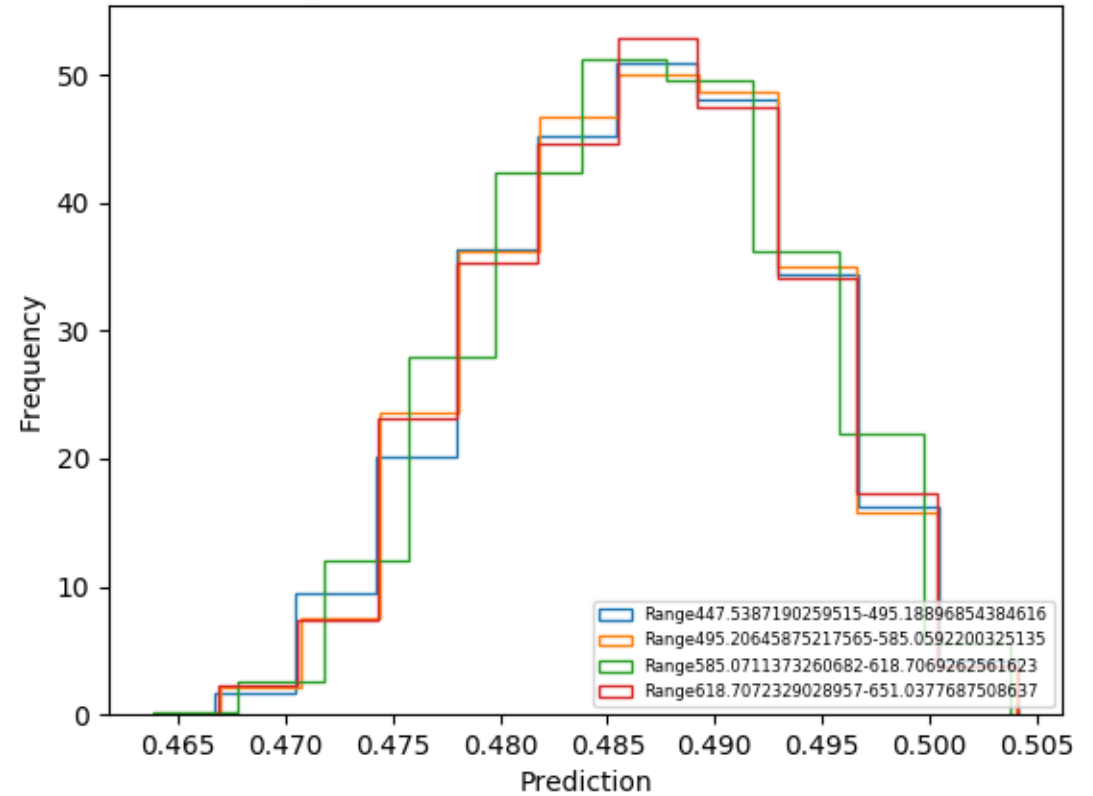


### 3. $D_s \rightarrow \eta (\rightarrow \mu^+ \mu^-) + \pi$ eta\_M 2016

Eta\_M Signal correlation for dseta 2016



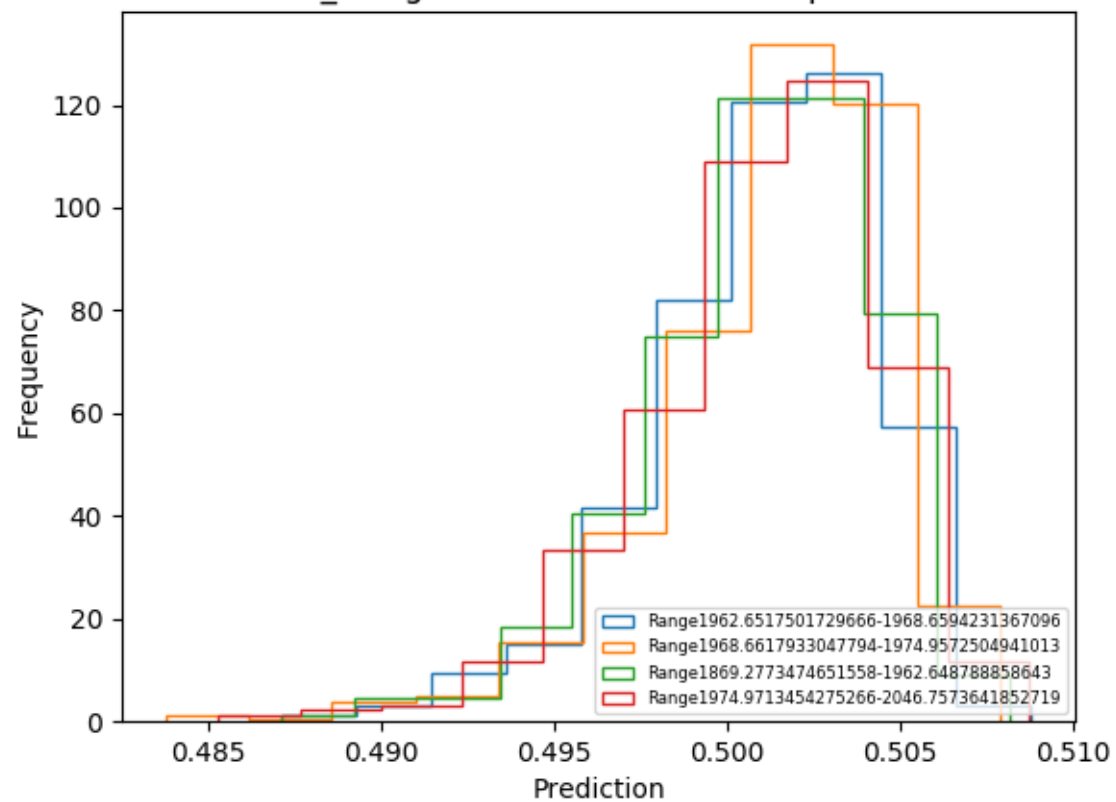
Eta\_M Background correlation dseta 2016



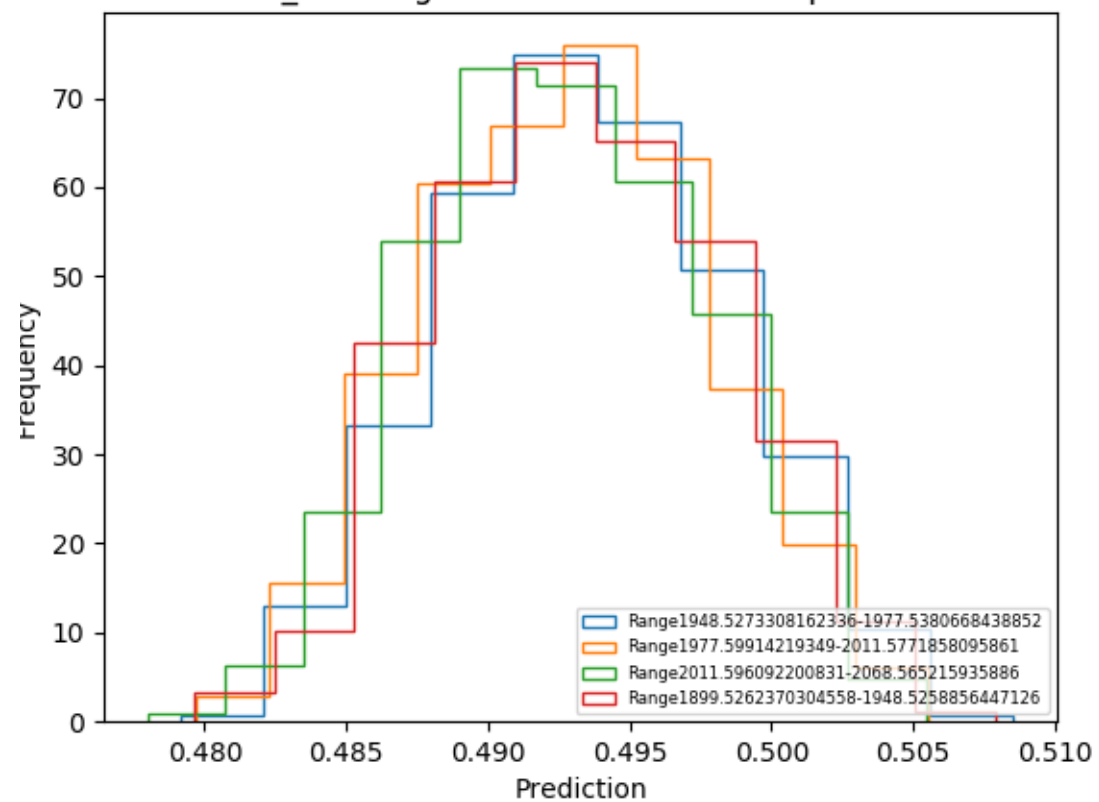


### 3. $Ds \rightarrow \eta' (\rightarrow \mu^+ \mu^-) + \pi$ Ds\_M 2017

Ds\_M Signal correlation for dsetap 2017



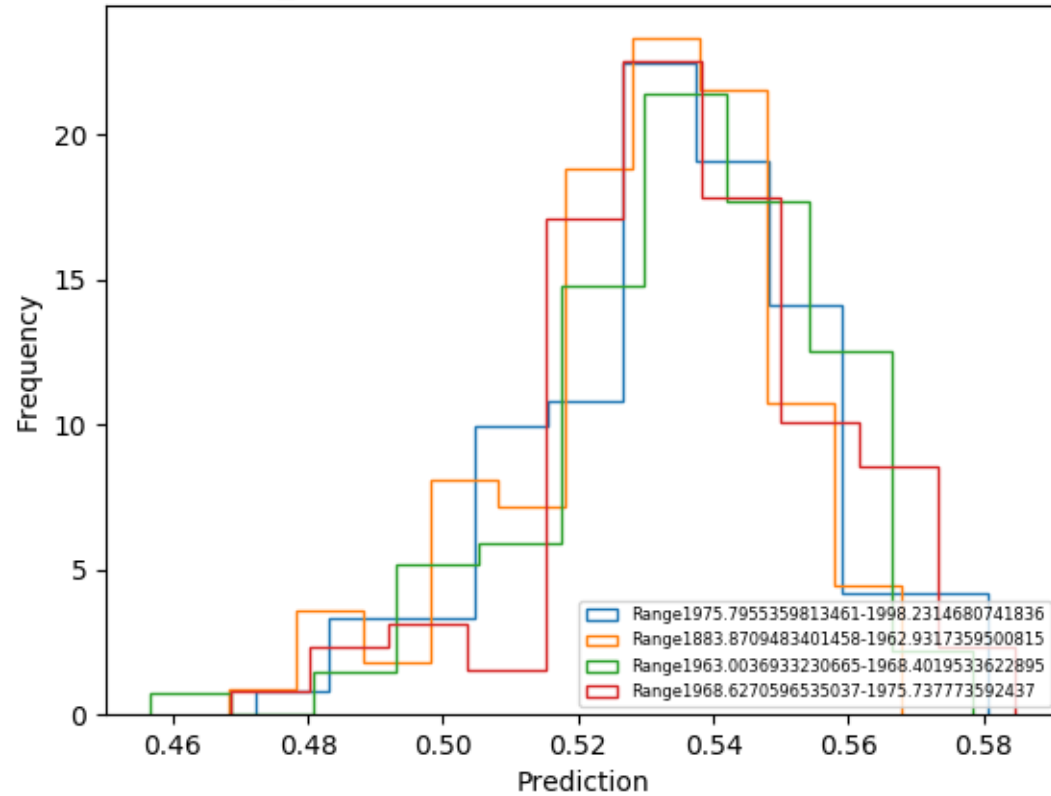
Ds\_M Background correlation dsetap 2017



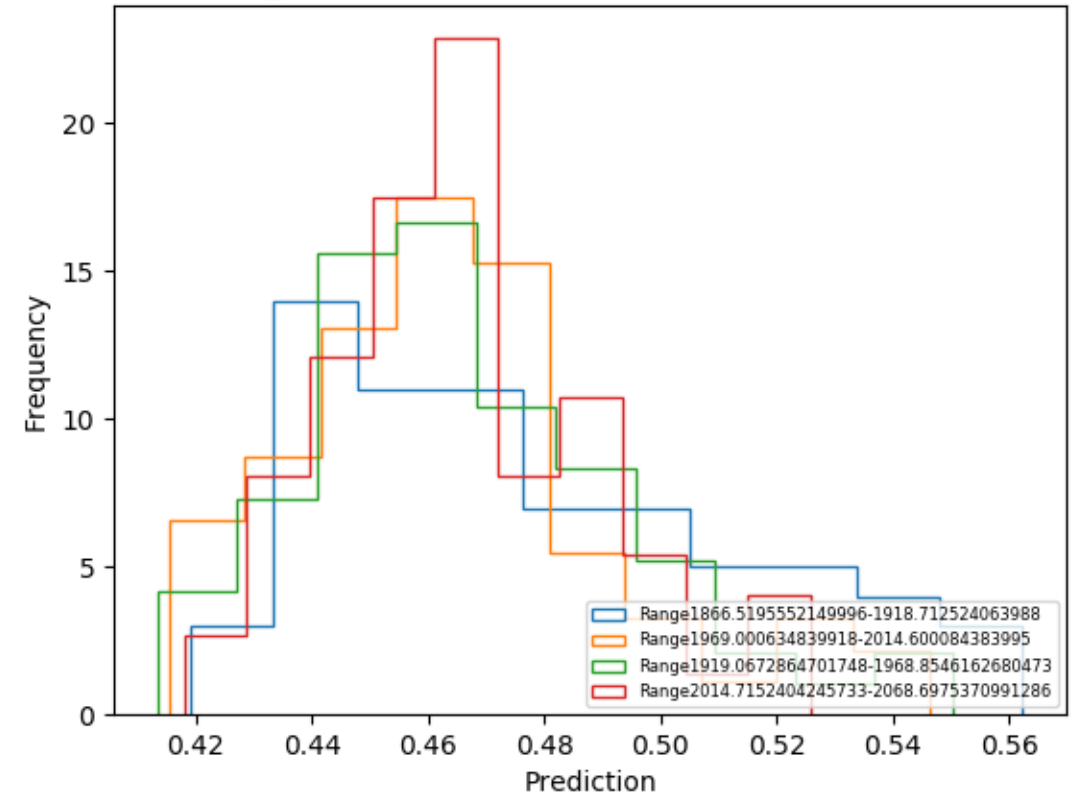


### 3. $Ds \rightarrow \eta' (\rightarrow \pi^+ \pi^- \mu^+ \mu^-) + \pi$ Ds\_M 2018

Ds\_M Signal correlation for mumupipi 2018



Ds\_M Background correlation mumupipi 2018



# Future work

---





As we showed before the future lines of work are:

- 1) The necessity of proving with different data that this conclusion is correct and not just a singular occasion in this specific data.
- 2) Proving why the results under the same specifications work better for our module in scikit-learning Python than our module in TMVA Python.
- 3) The development of Machine Learning model to solve problems related to the industry.
  - 1) So far we started a project related with a company in the textile sector to detect and solve failures in the system.



TRIPLE**ALPHA**  
innovation