

Machine Learning applications in high energy physics

Viraj Thakkar

MSc Thesis Presentation

Supervised by Dr. Bedangadas Mohanty

NISER, Bhubaneswar

Motivation

- The aim of this project is to use **Machine Learning** techniques for the purpose of classification of signal and background events of resonance particle K^{*0} and **improve the significance** of the signal as compared to traditional approach.

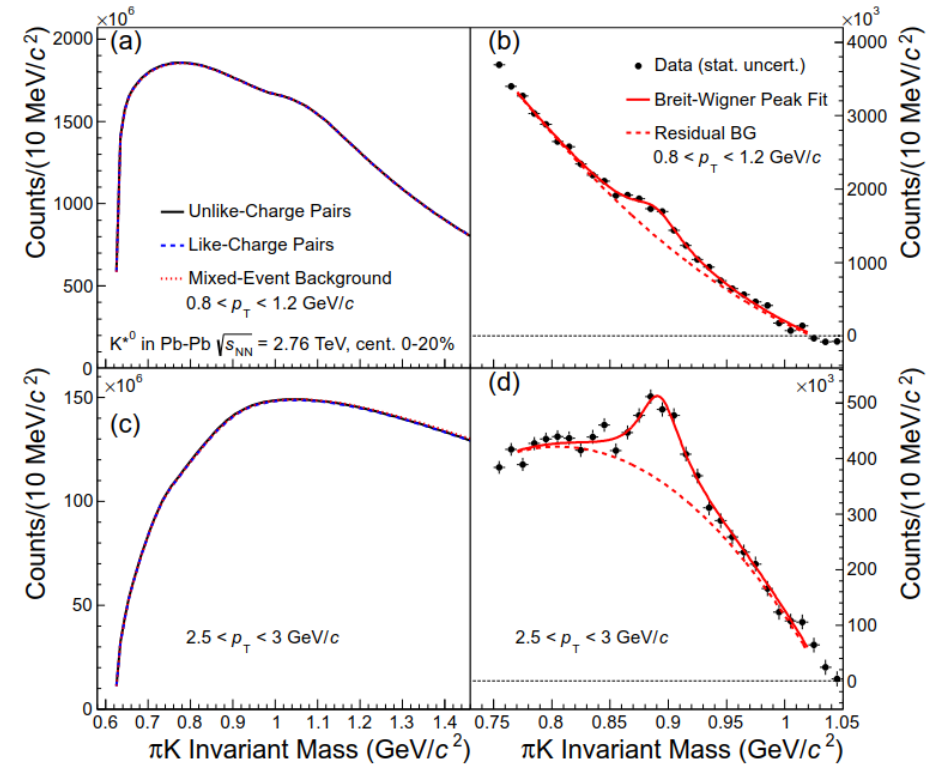


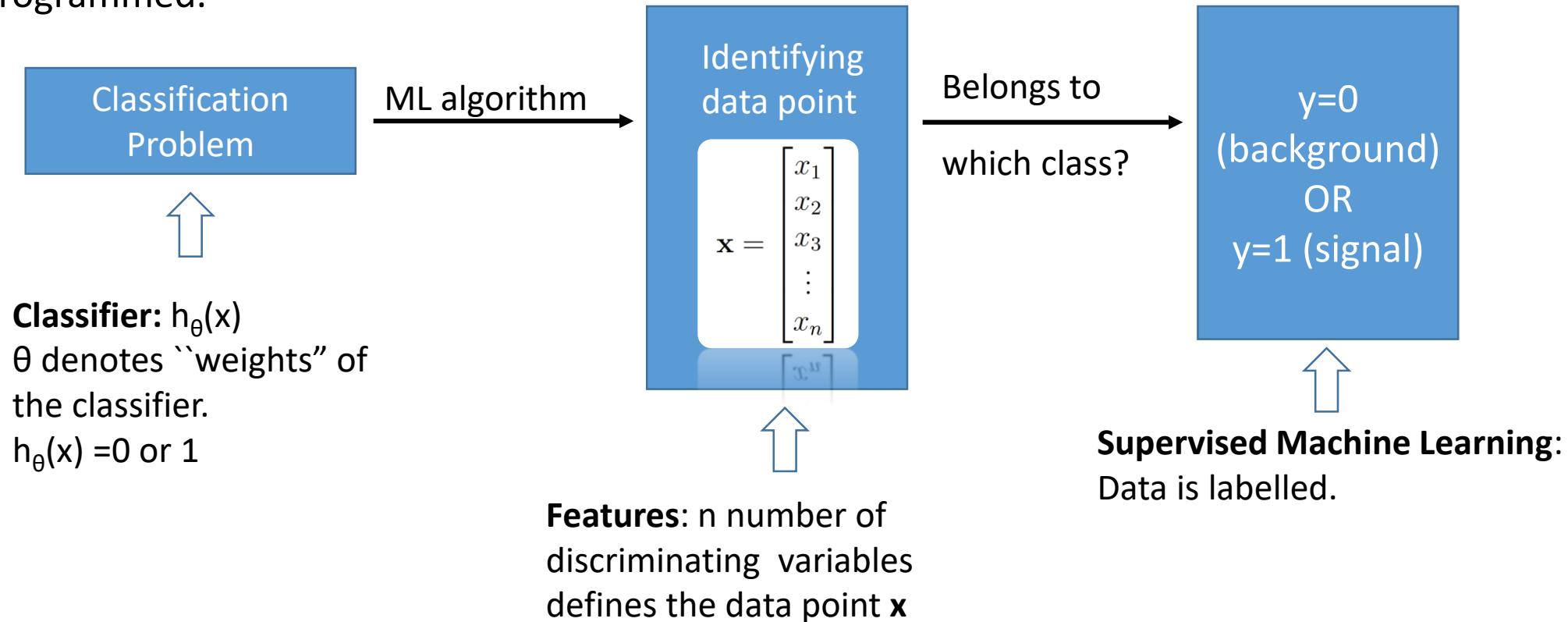
Fig: Large combinatorial background (left), signal is rarely visible

Reference: B.Abelev et al, ALICE

Collaboration: Phys. Rev. C 91 (2015) 024609

Machine Learning

➤ Machine Learning(ML) is a field of Artificial Intelligence(AI) that uses statistical techniques to provide computer systems the ability to “learn” from data, without being explicitly programmed.



The Machine Learning process



Training

→The classifier uses data points \mathbf{x} and learns about the features

→Classifier learns the data in the form of weights θ

→Uses the label $y=0$ or 1 for the data point \mathbf{x}

Testing

→The classifier makes predictions on untrained data points \mathbf{x}

→Compare the classifier predicted output $h_{\theta}(\mathbf{x})$ with the actual label y of the data \mathbf{x}

→Does not use the labels $y=0$ or 1 for the data point \mathbf{x} while testing.

Application

→The classifier makes predictions on new unknown data points \mathbf{x}

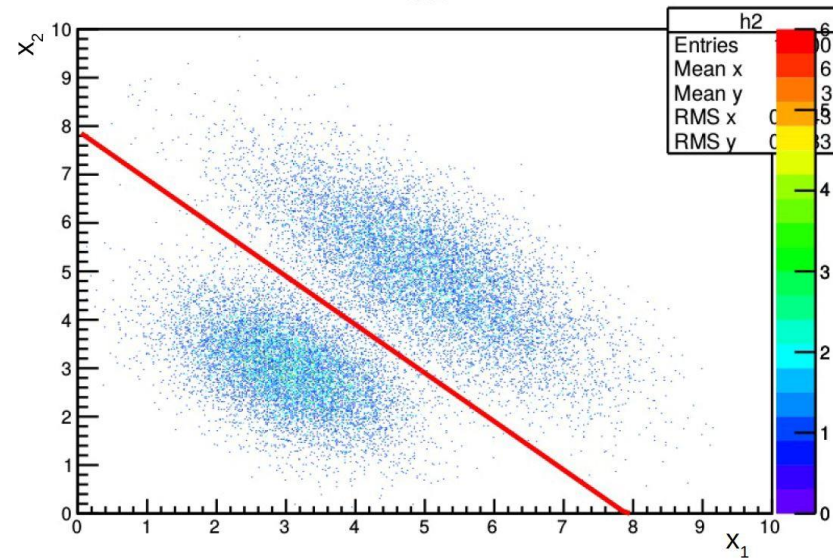
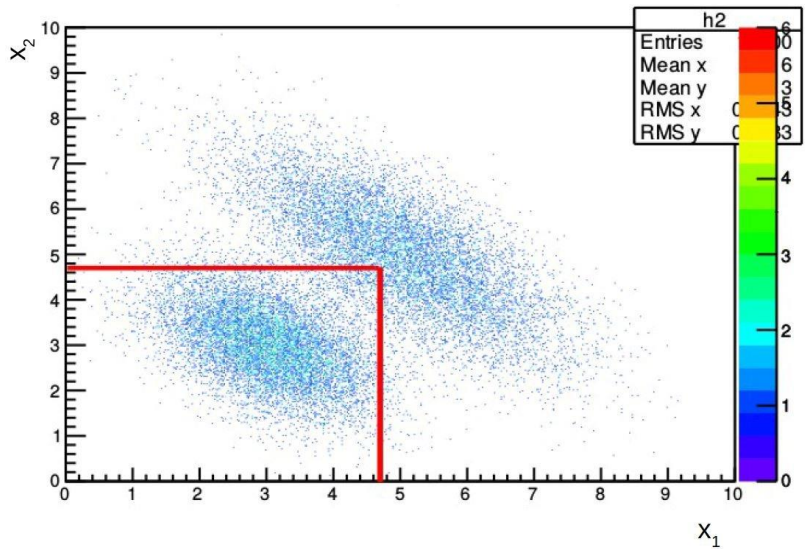
→Predictions are made using the weights θ which the classifier learned during training

→Labels $y=0$ or 1 are not known. Classifier predicts it.

Toolkit For Multivariate Analysis (TMVA)

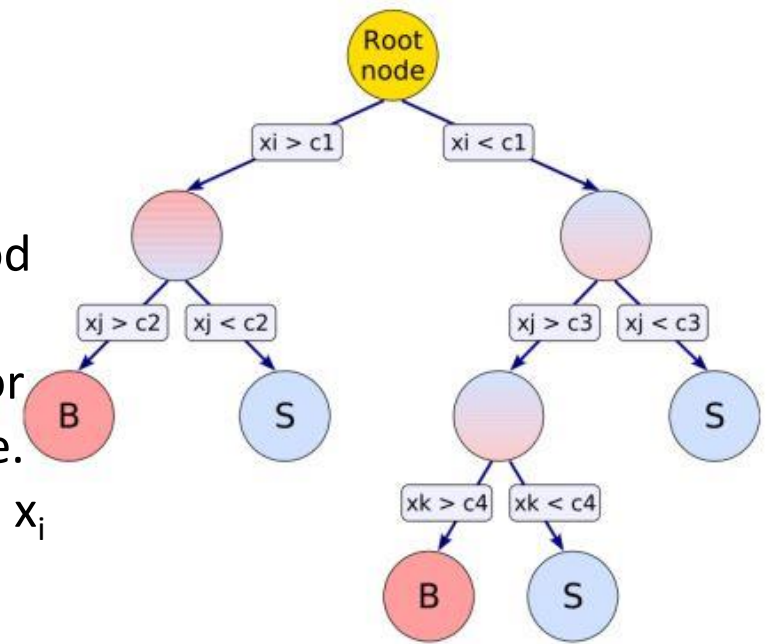
- TMVA provides a framework for multivariate analysis (ML classification algorithms) implemented in ROOT
- We have used Boosted Decision Tree(BDT) as our ML classification algorithm.

Why we need Multivariate Analysis? : Better separation



Decision Trees

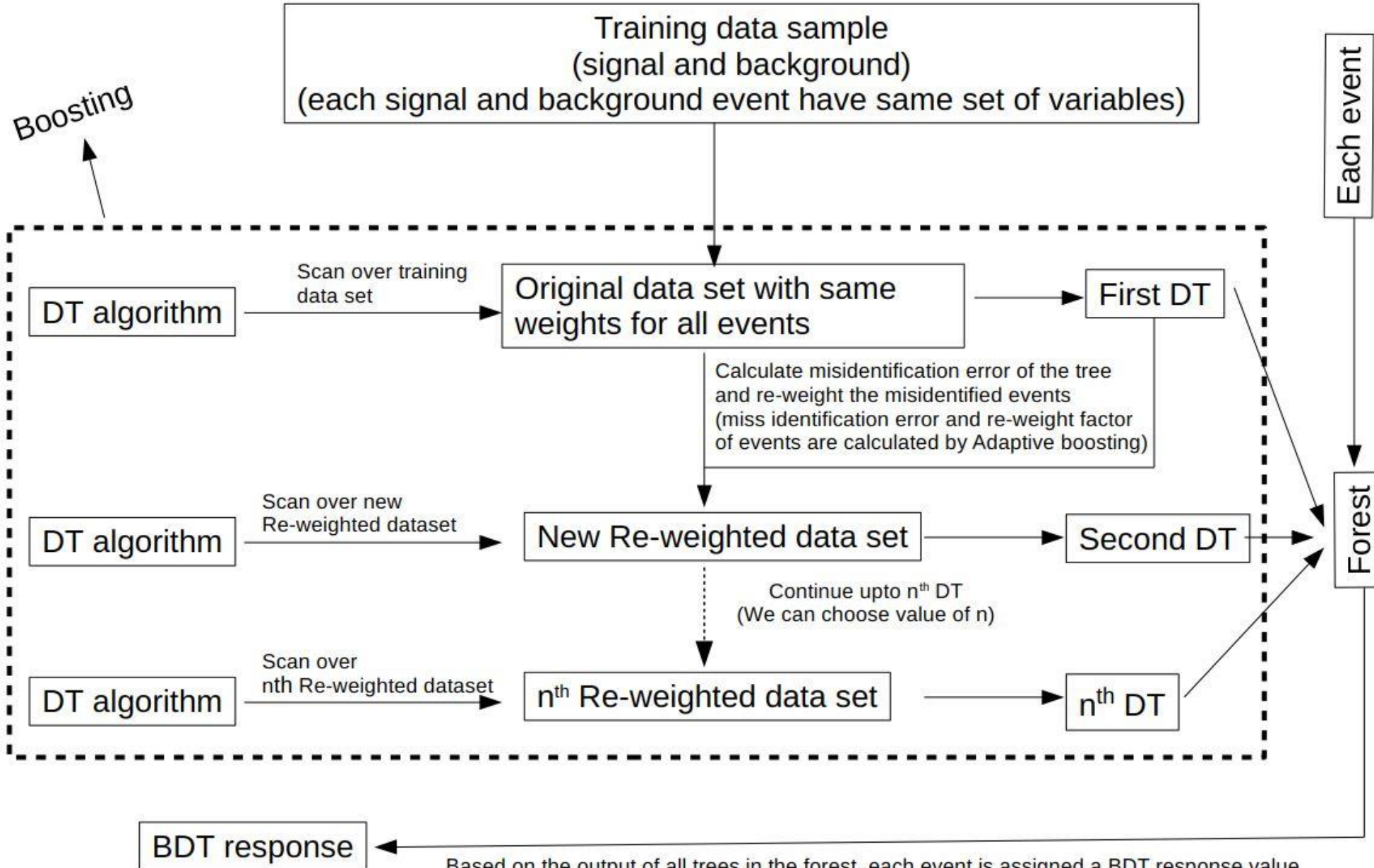
- A decision tree(DT) is a tree-like structure that uses a branching method to illustrate possible outcome of a decision.
- An event is either classified as signal or background by either passing or not passing a condition (cut) on a specific node until a decision is made.
- A sequence of binary splits using the discriminating variables(features) x_i is applied to the data.
- The split is such that we get the best separation between signal and background.



Boosting

- Boosted Decision Trees (BDT) is a prediction model which uses an ensemble of “weak classifiers” (decision trees).
- It is a model designed to make fewer and fewer mistakes as more trees are added to it.
- BDT is an algorithm which combines forests of DTs, each weighted according to their importance.

Boosting



Based on the output of all trees in the forest, each event is assigned a BDT response value ranging from -1 to +1. Background-like events will have a BDT response value shifted towards -1, and signal-like events will have a response value shifted towards +1

Data Set

We have used pp collisions data at $\sqrt{s} = 7$ TeV from ALICE detector at LHC

Reconstructed Monte Carlo

System: p+p

Period: LHC10f6a

Event generator: PYTHIA

Events: ~10M

Data

System: p+p

Period: LHC10d

Events: ~10M

Event Selection: $|V_z| < 10$ cm.

Track Selection: ITS-TPC 2010 cuts

Reference: B. Abelev et al ALICE Collaboration: Eur. Phys. J. C72 (2012) 2183

Preparation of data for training, testing and application



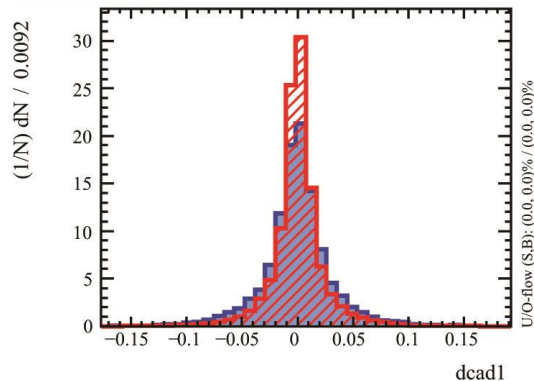
- $K^{*0} :=$ Mother Particle
- $K^+ :=$ Daughter Particle 1
- $\pi^- :=$ Daughter Particle 2

- **Signal** candidates are formed by K^+ and π^- pairs which are decayed from K^{*0} . Generated level MC PID is used to select true K^+ and π^-
- **Background** candidates are formed by same sign K and π like pairs which are selected by using energy loss information in TPC detector.
- **Unlike** pairs are all combinations of K^+ and π^- pairs which constitutes signal plus background events.

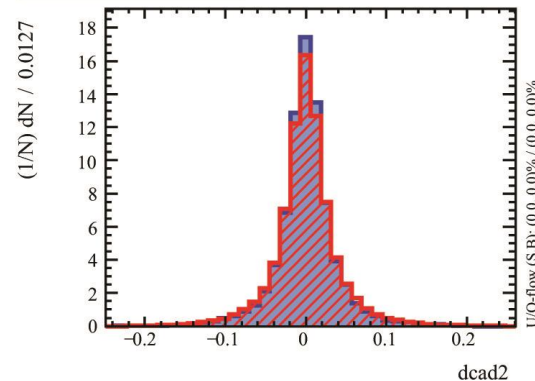
Input features: MC (PYTHIA)

$0 < p_T(\text{mother}) < 0.8 \text{ GeV}/c$

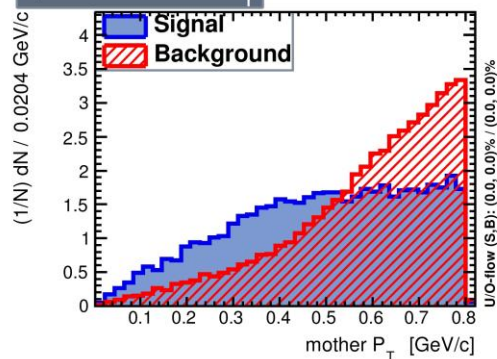
Input variable: dcad1



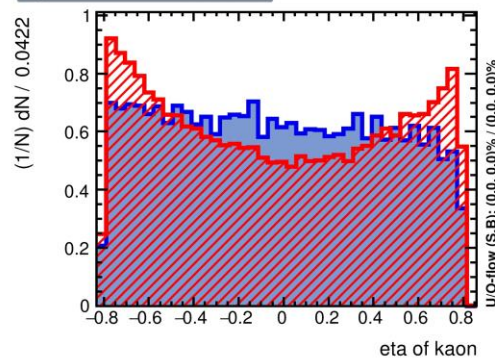
Input variable: dcad2



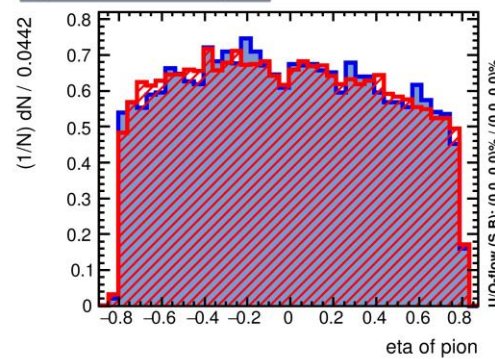
Input variable: mother P_T



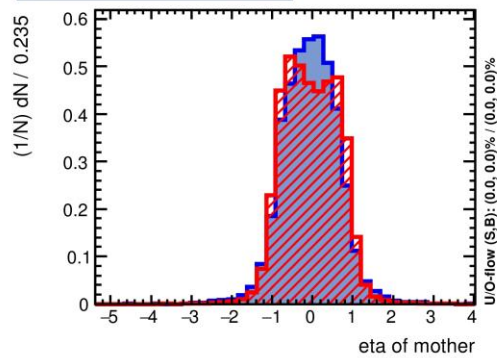
Input variable: eta of kaon



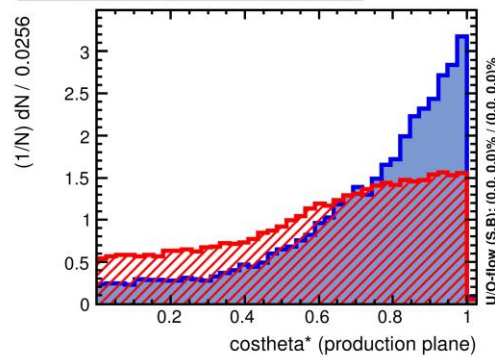
Input variable: eta of pion



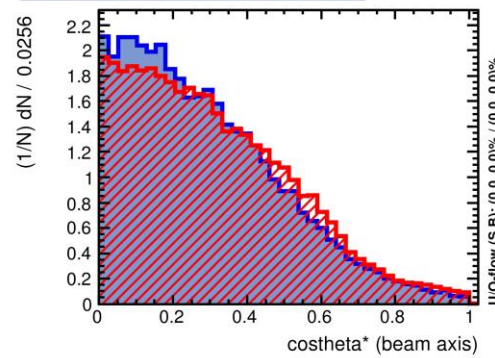
Input variable: eta of mother



Input variable: costheta* (production plane)

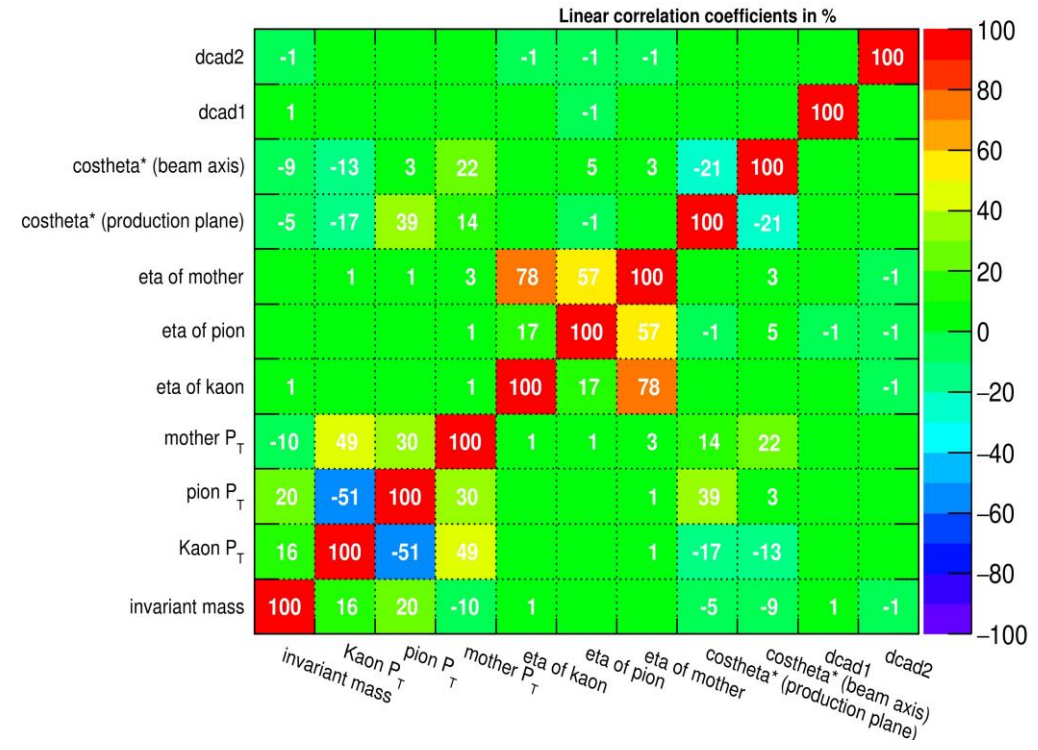


Input variable: costheta* (beam axis)

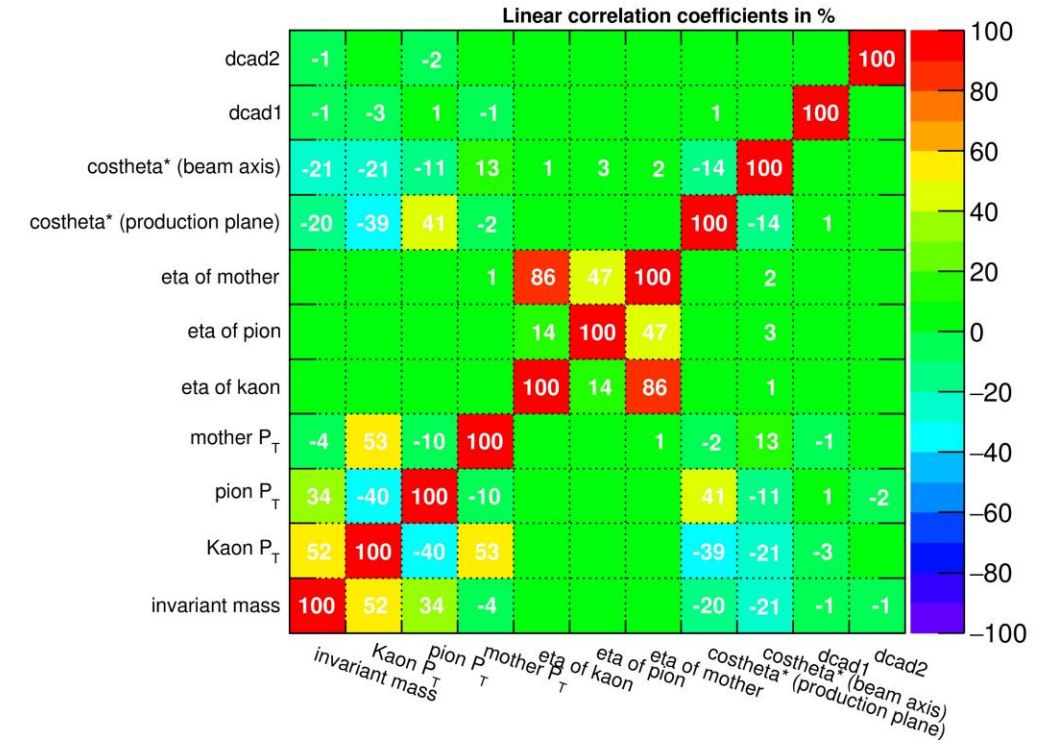


Correlation coefficient percentage

Correlation Matrix (signal)



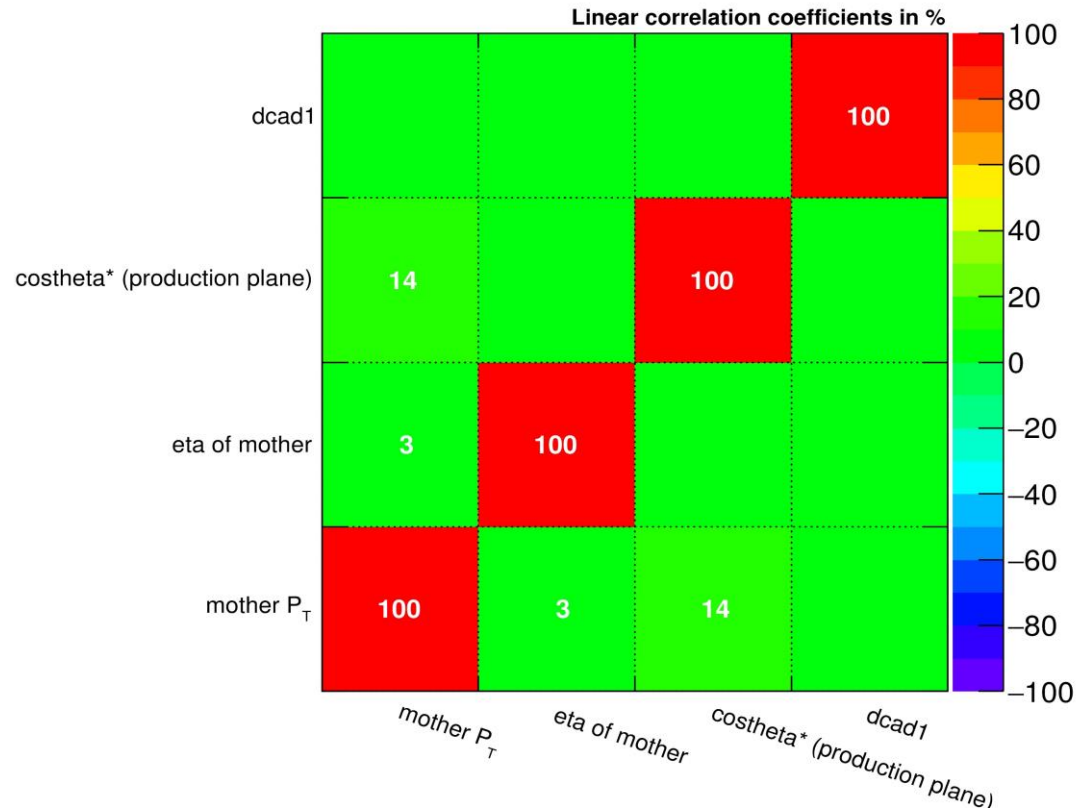
Correlation Matrix (background)



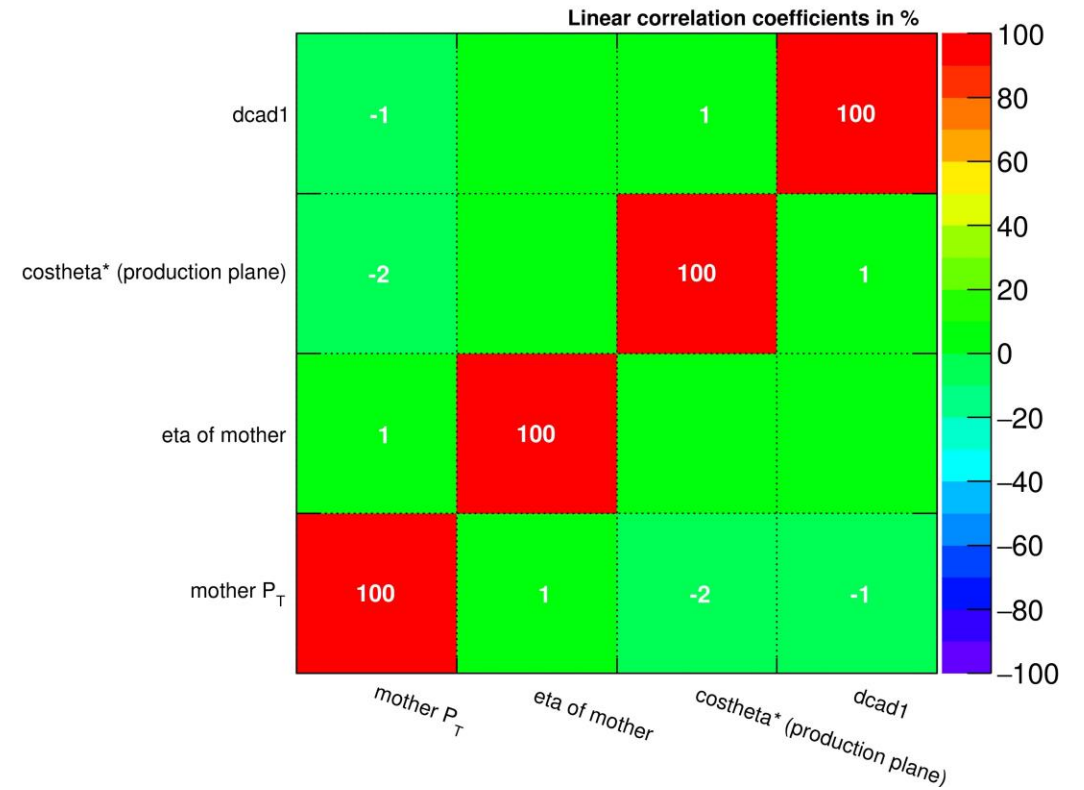
Input features used for training phase

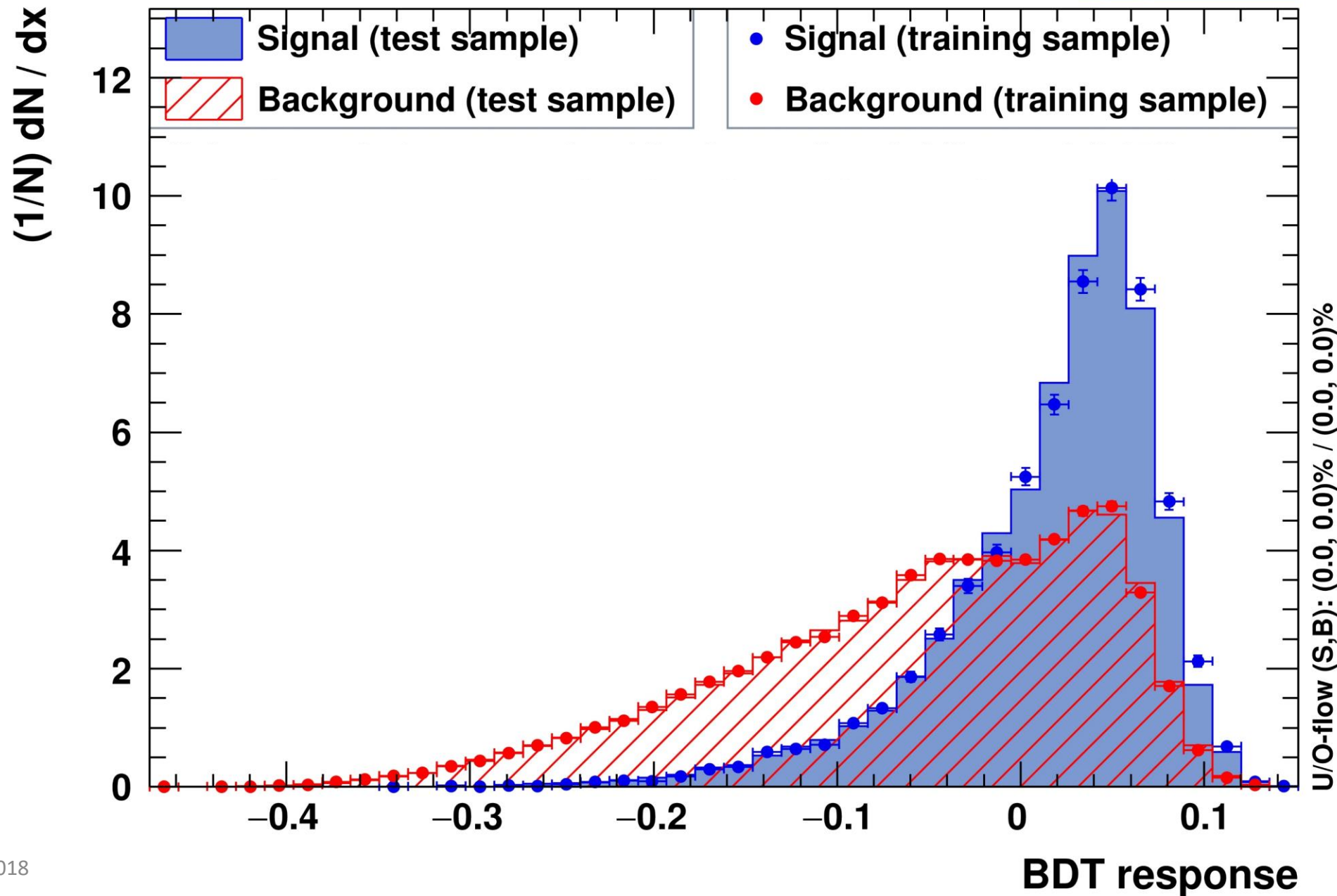
Selected features have less correlation with invariant mass

Correlation Matrix (signal)

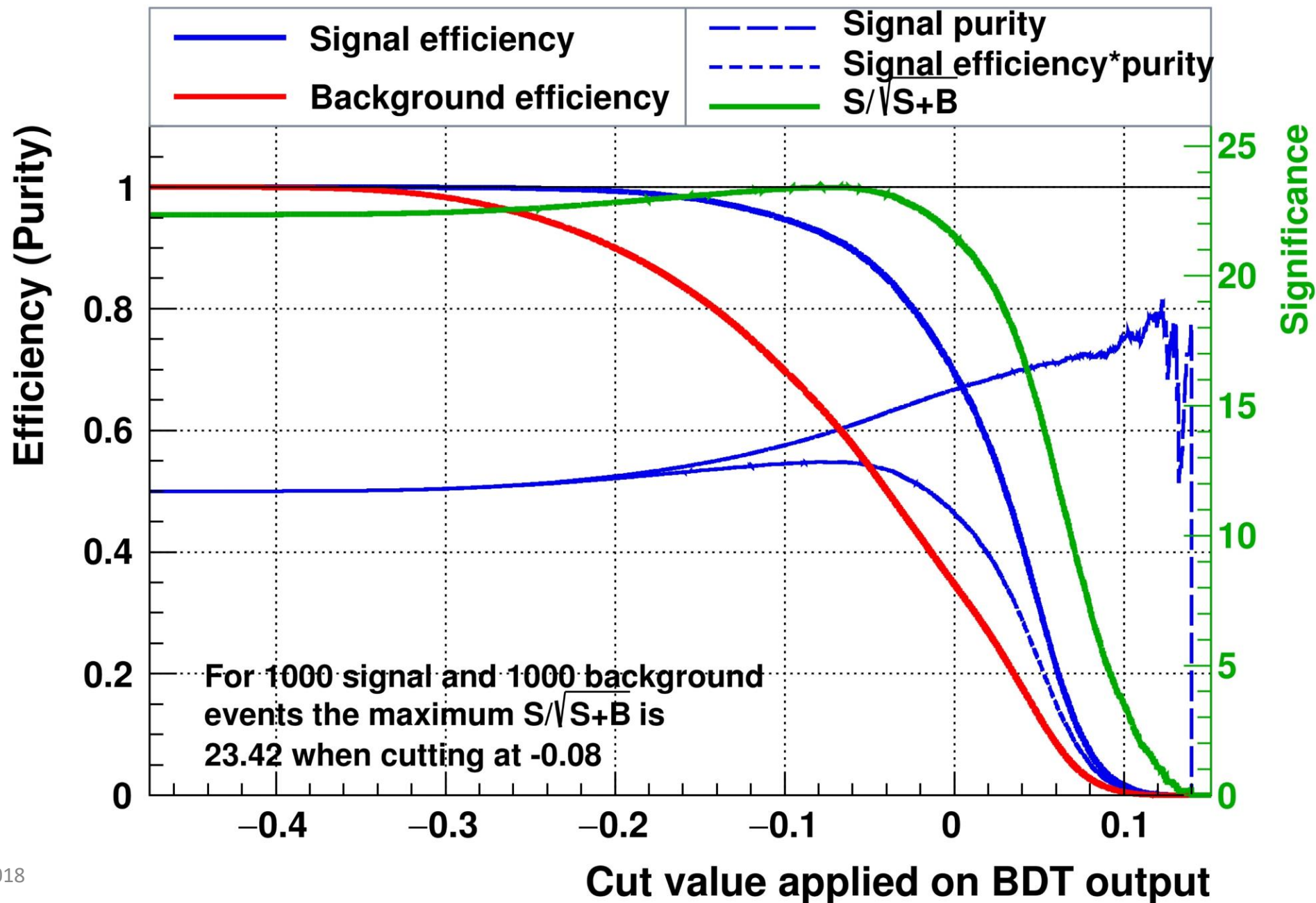


Correlation Matrix (background)

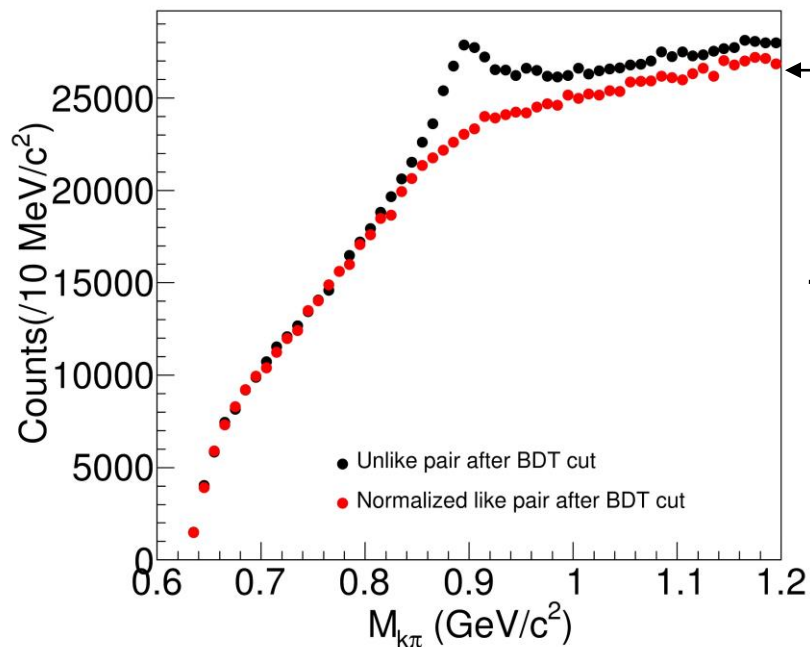
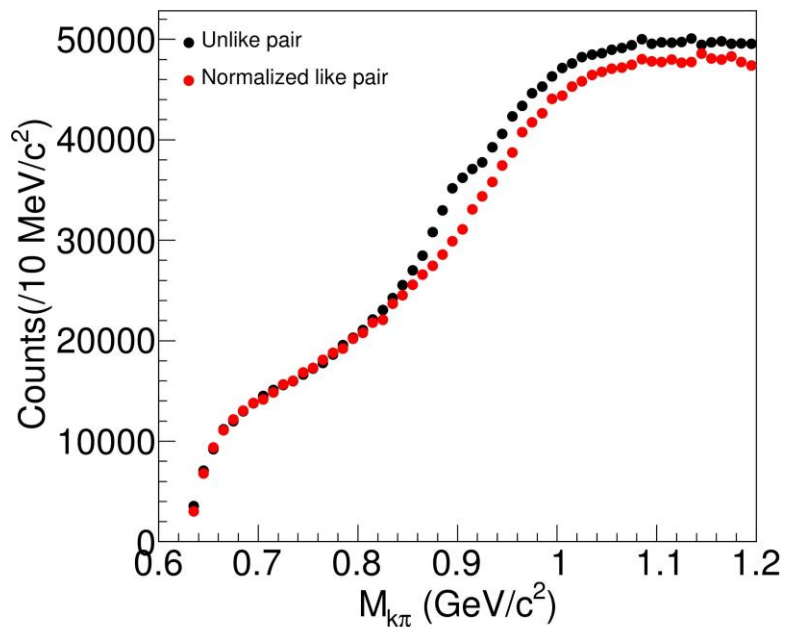




Cut efficiencies and optimal cut value



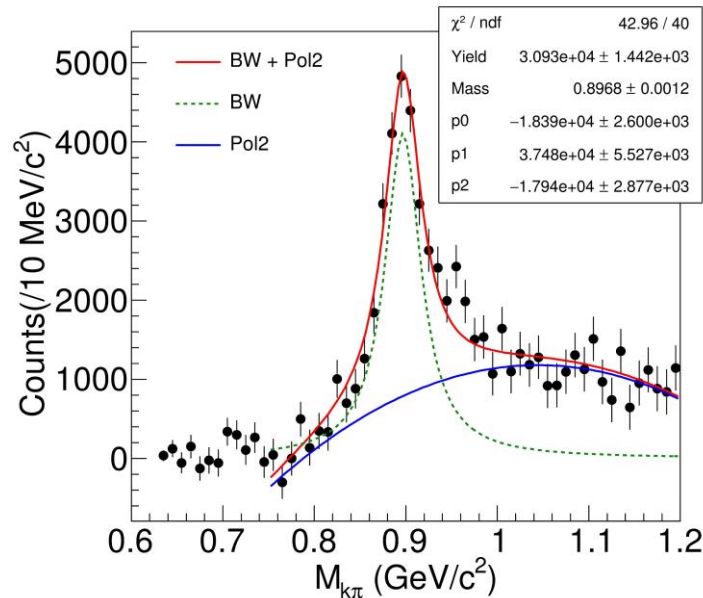
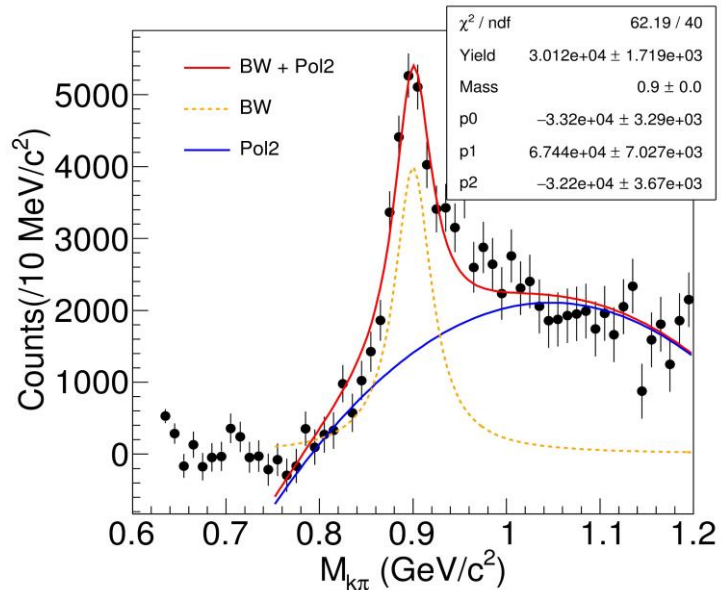
Invariant mass distribution: MC (PYTHIA)



Signal is more visible!

ML approach with BDT

Traditional approach



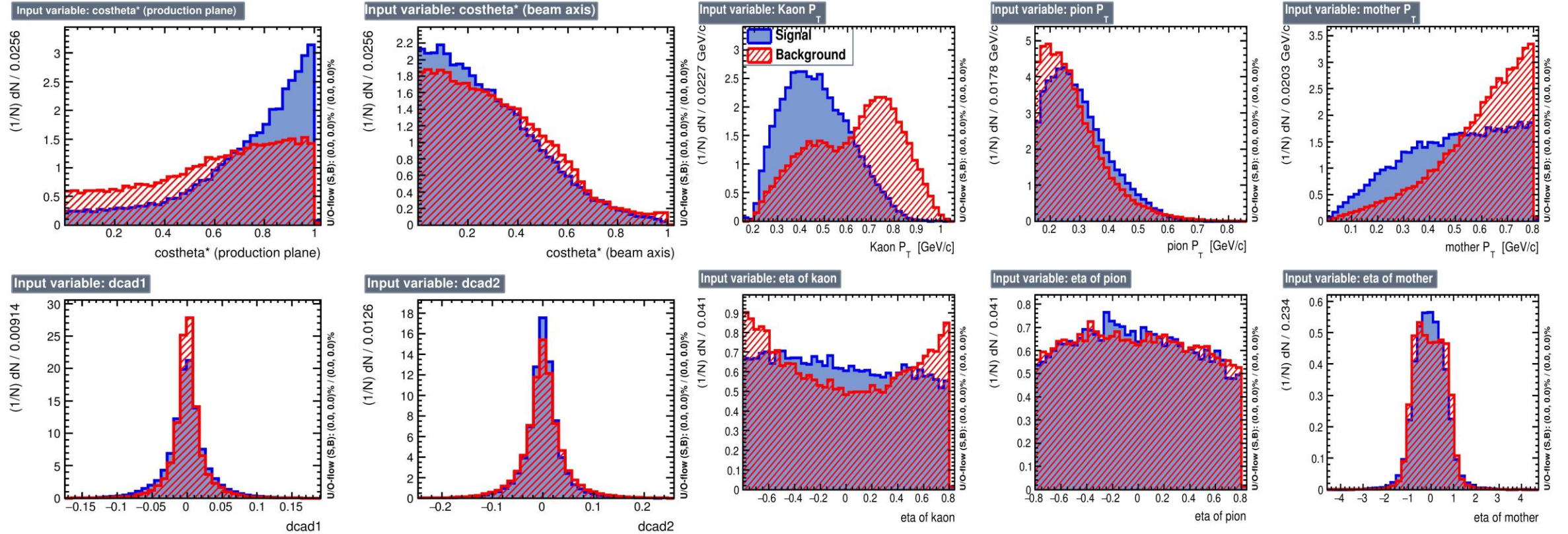
Comparison of signal obtained by Monte-Carlo

Method	χ^2/ndf	S	$S + B$	$S/\sqrt{S + B}$	Input Yield
Machine Learning	42.96/40	$\frac{30931.6}{0.92}$ = 33621.3	906953	32.48	34253
Traditional Approach	62.19/40	30115	1355500	25.87	34253

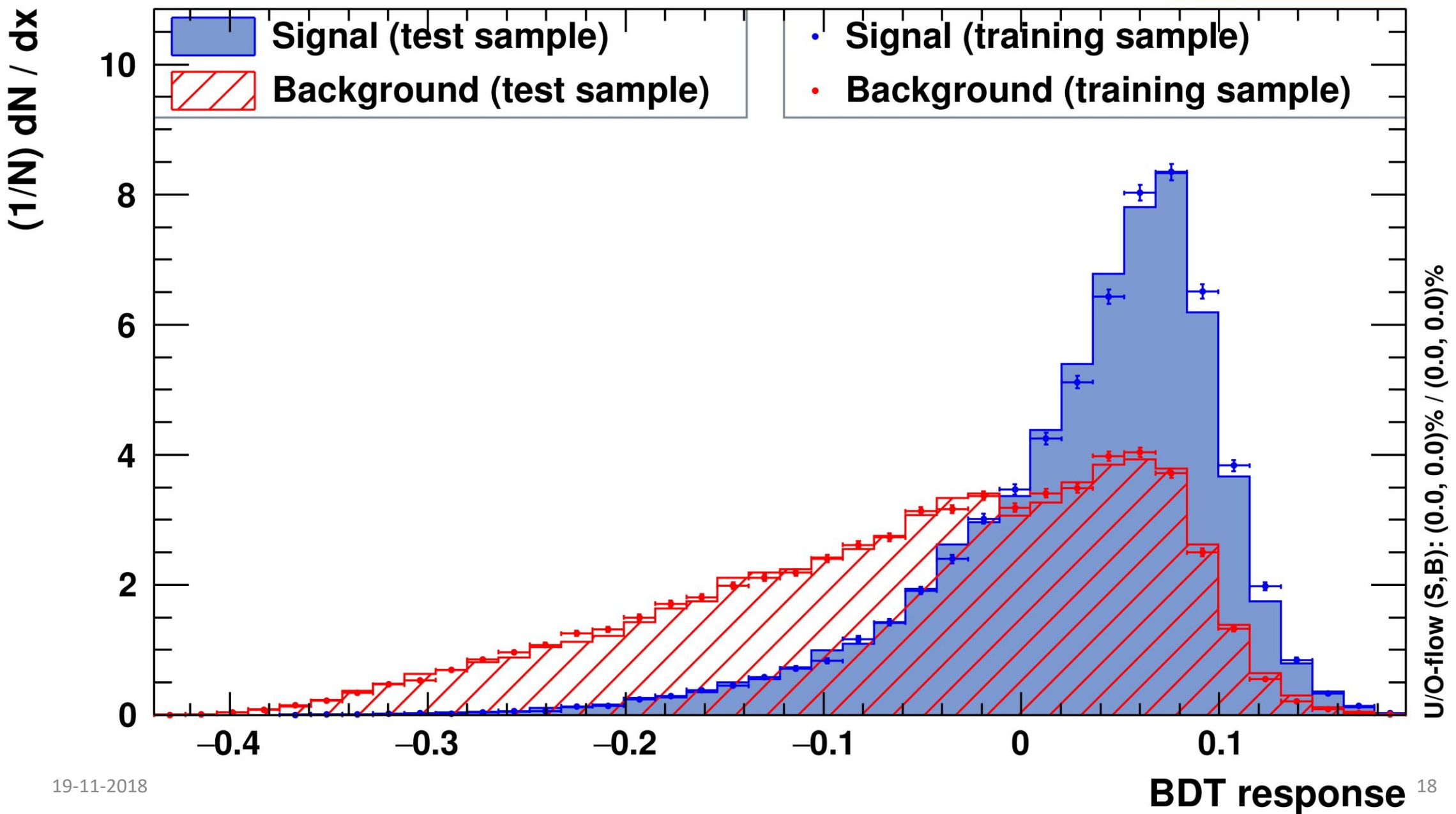
Extraction yield, significance and χ^2/ndf of the fit suggests a better extraction of signal by using Machine Learning via BDT for MC data.

Input Features: ALICE data

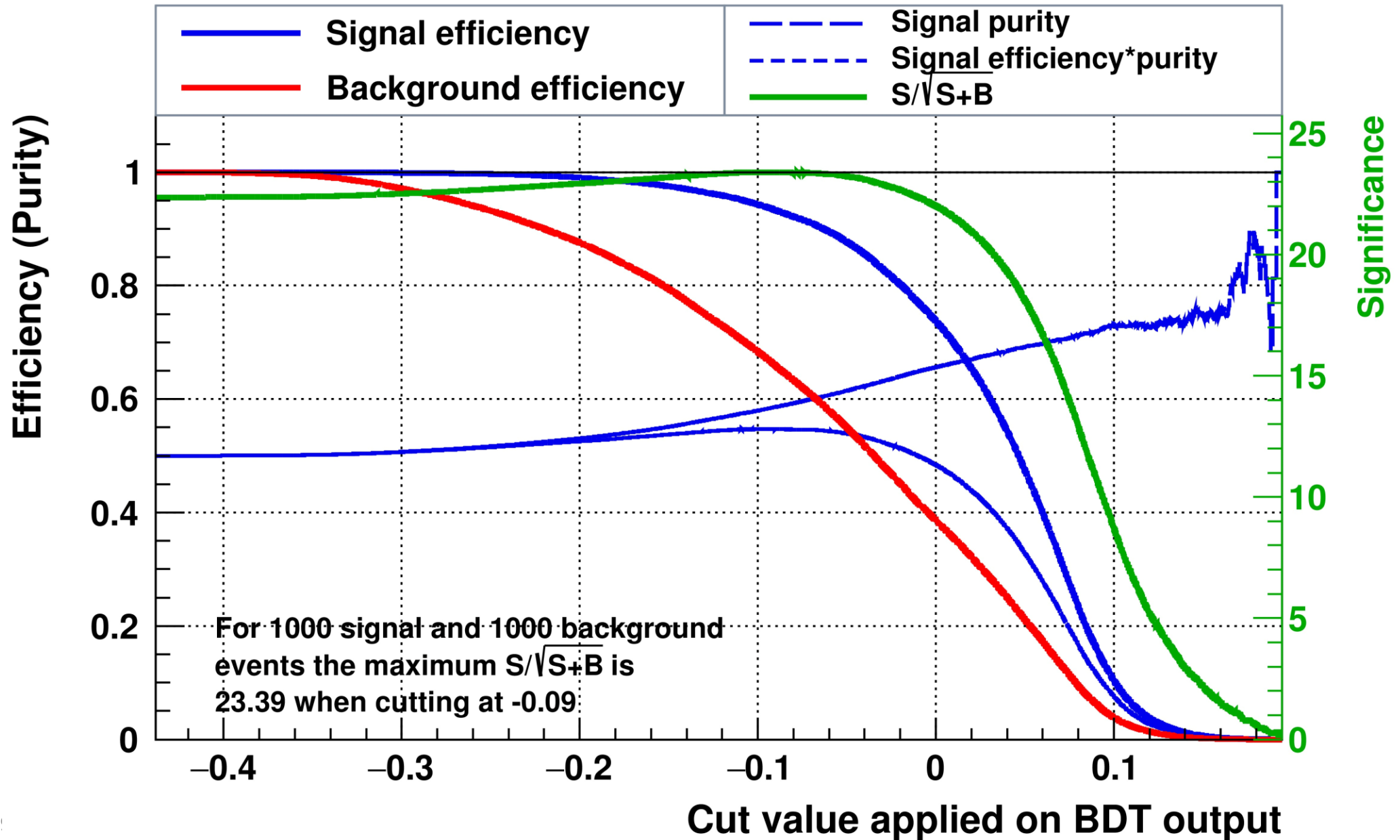
$0 < p_T(\text{mother}) < 0.8 \text{ GeV}/c$



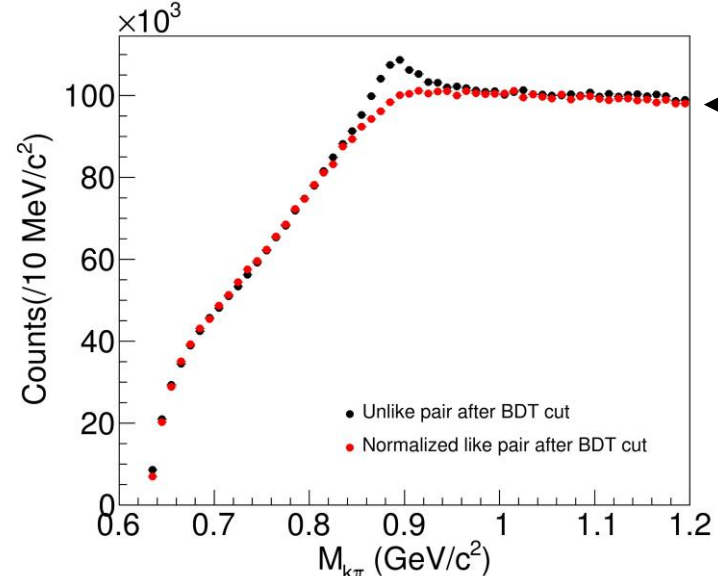
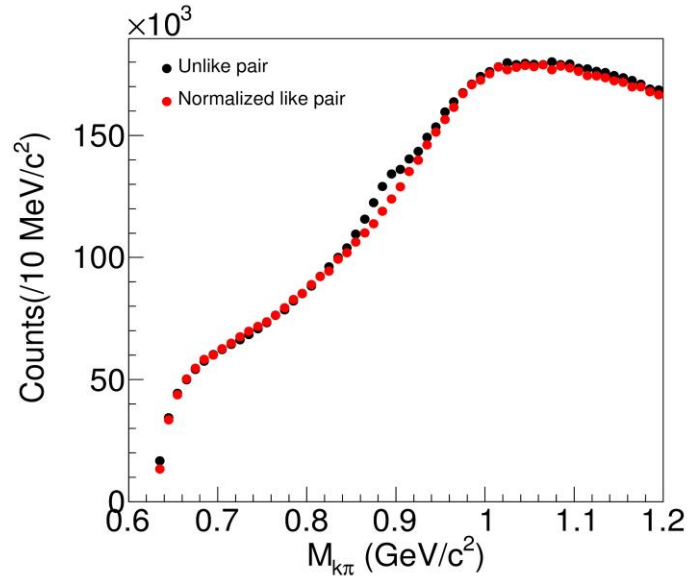
TMVA overtraining check for classifier: BDT



Cut efficiencies and optimal cut value



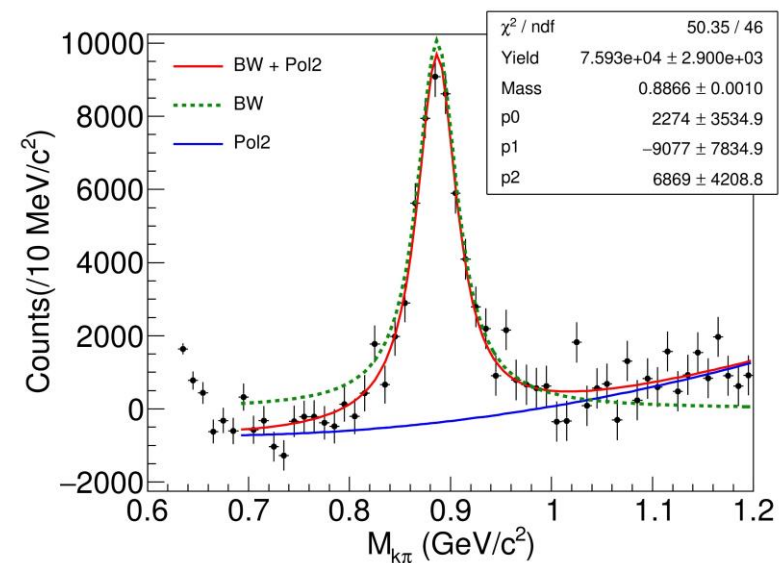
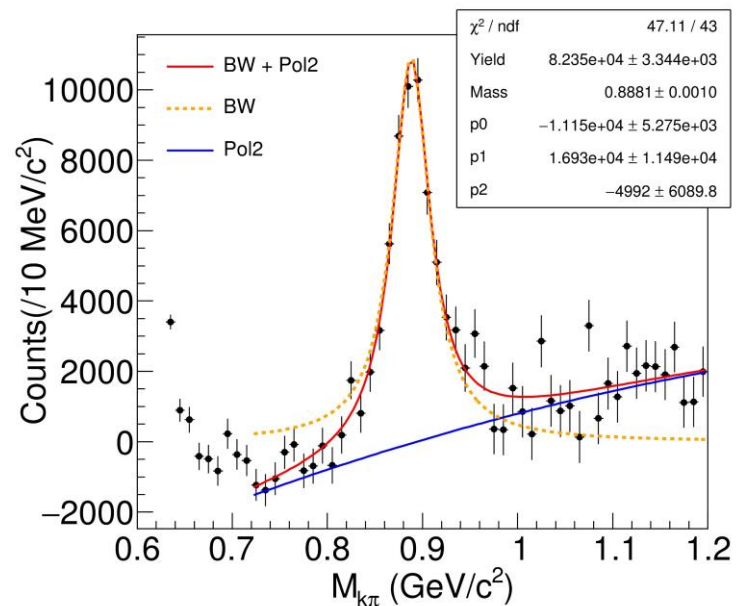
Invariant mass distribution: ALICE data



Signal is more visible!

Traditional approach

ML approach with BDT



Comparison of signal obtained from ALICE data

Method	χ^2/ndf	S	$S + B$	$S/\sqrt{S + B}$
Machine Learning	50.35/46	$\frac{75934.8}{0.93}$ $= 81650.3$	3626890	38.54
Traditional Approach	47.11/43	82346.8	5227420	34.82

Extraction yield, significance and χ^2/ndf of the fit suggests a better extraction of signal by using Machine Learning via BDT for real data too.

Summary

In this analysis, we have extracted K^{*0} signal by two methods

1) Traditional invariant mass technique 2) ML classification by BDT

In both cases of Monte-Carlo and real data, we obtained a better significance by machine learning method as compared to the traditional method.

Outlook

- In future, we will extend this analysis for resonances which undergo 3-body decay.
- We will implement this analysis on Pb-Pb nuclear collisions data
- Machine Learning can also be applied to extract signal from rare resonances such as $K^*(1410)$, $K^*(1680)$ and $\Xi(1820)$.

Backup

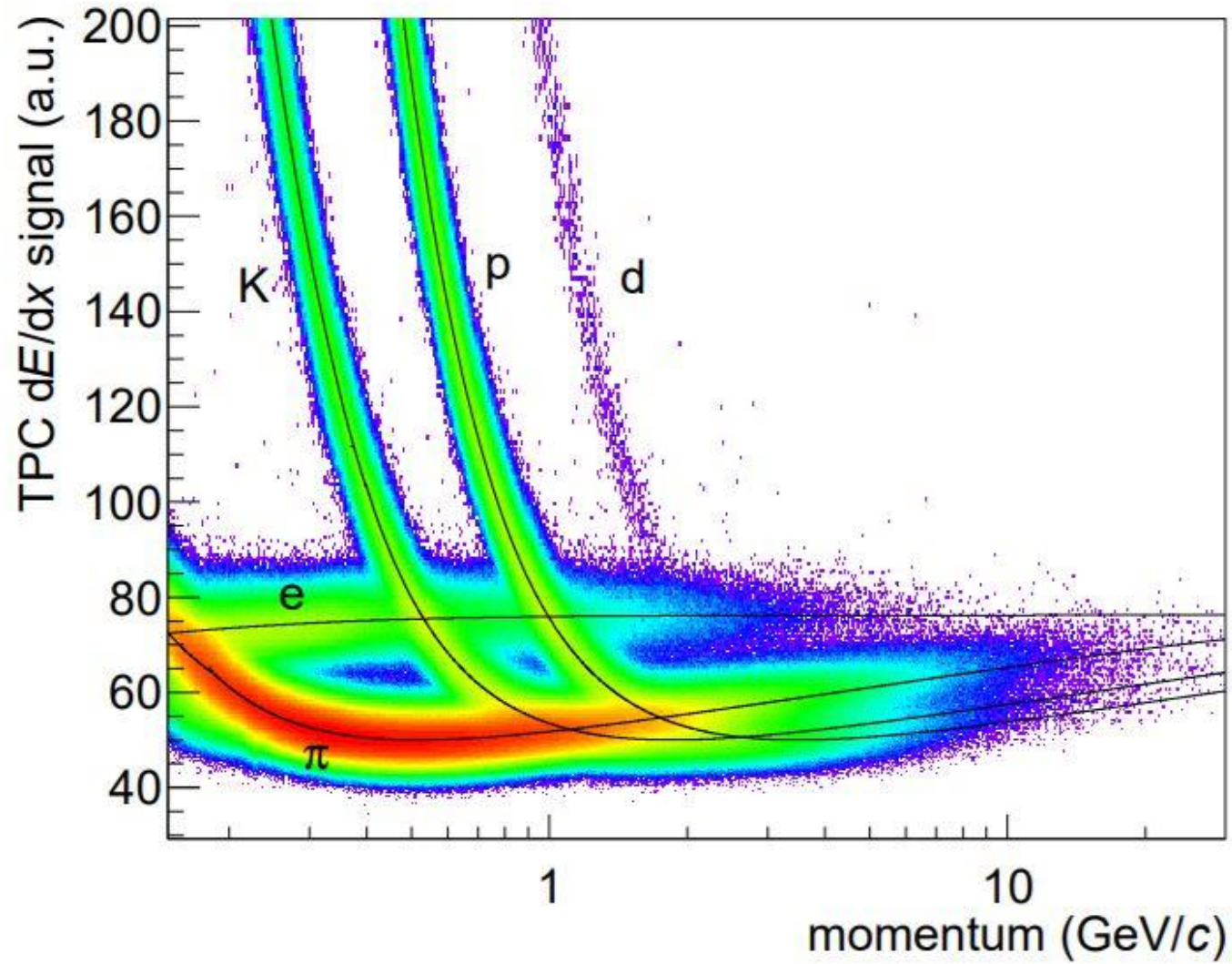
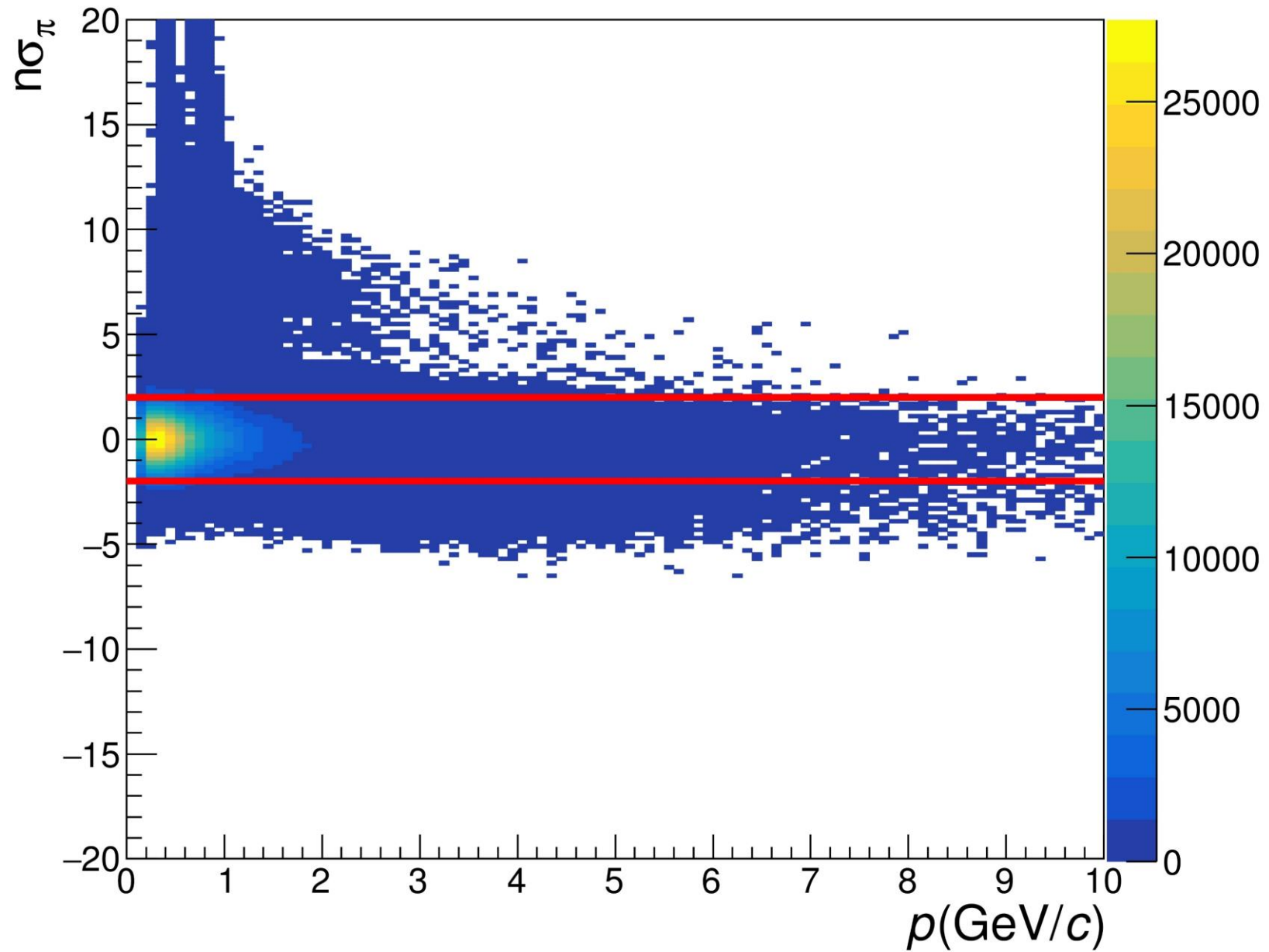


Fig. 1: (Colour online) Specific ionization energy loss dE/dx vs. momentum for tracks measured with the ALICE TPC. The solid lines are parametrizations of the Bethe-Bloch function [23].

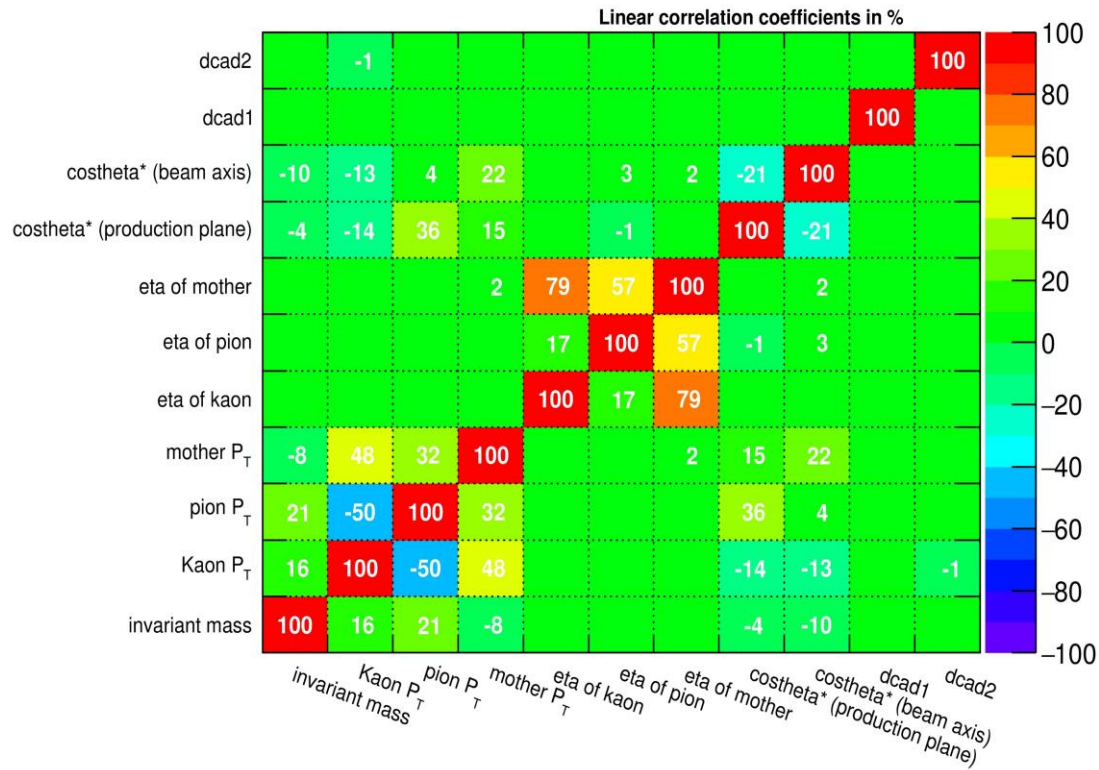


Standard Track Selection Cuts

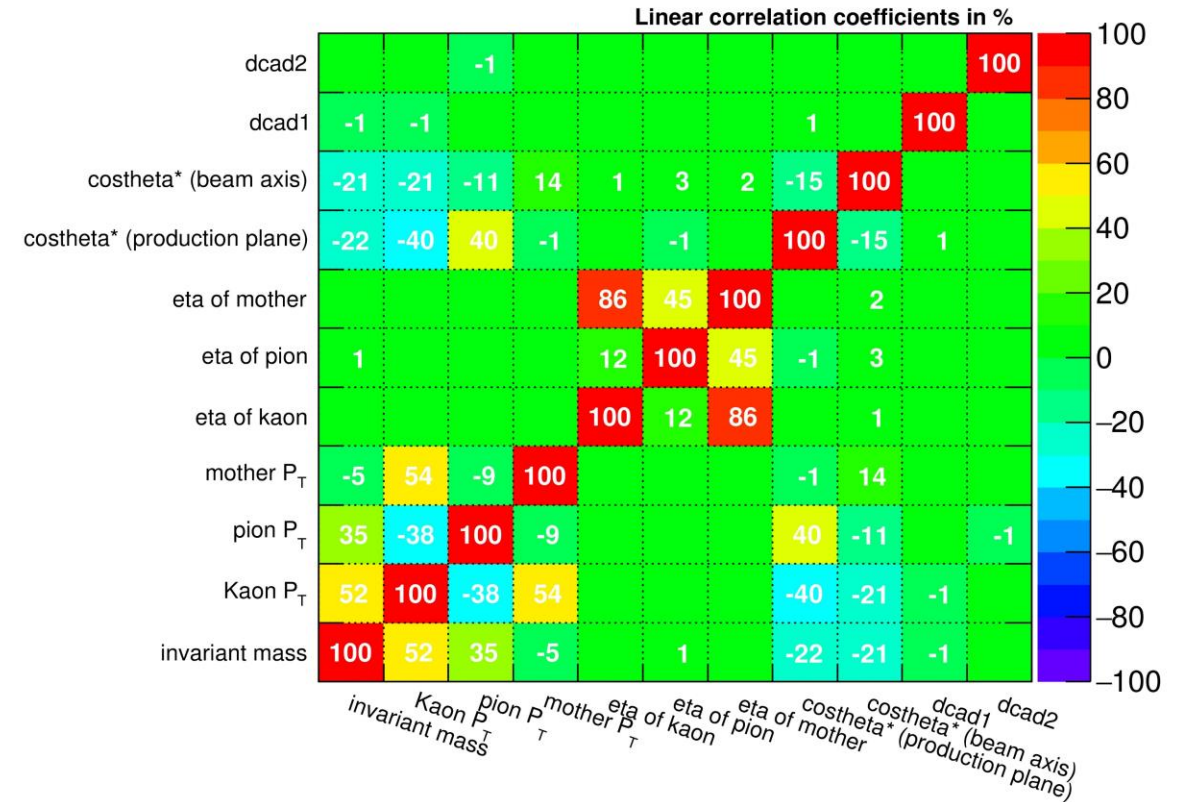
1. $p_T > 0.15 \text{ GeV}/c$
2. $-0.8 < \eta < 0.8$
3. Reject kink daughters
4. Ratio of crossed rows over findable cluster > 0.8
5. Minimum(Maximum) number of rows crossed in TPC is 70(159).
6. TPC $\chi^2/\text{clusters} < 4.0$
7. ITS $\chi^2/\text{clusters} < 36.0$
8. $(DCA)_r(p_T) < 0.0182 + 0.0350p_T^{-1.0} \text{ cm}$
9. $|DCA|_z < 2 \text{ cm}$
10. $|y_{pair}| < 0.5$

Correlation coefficient percentage

Correlation Matrix (signal)



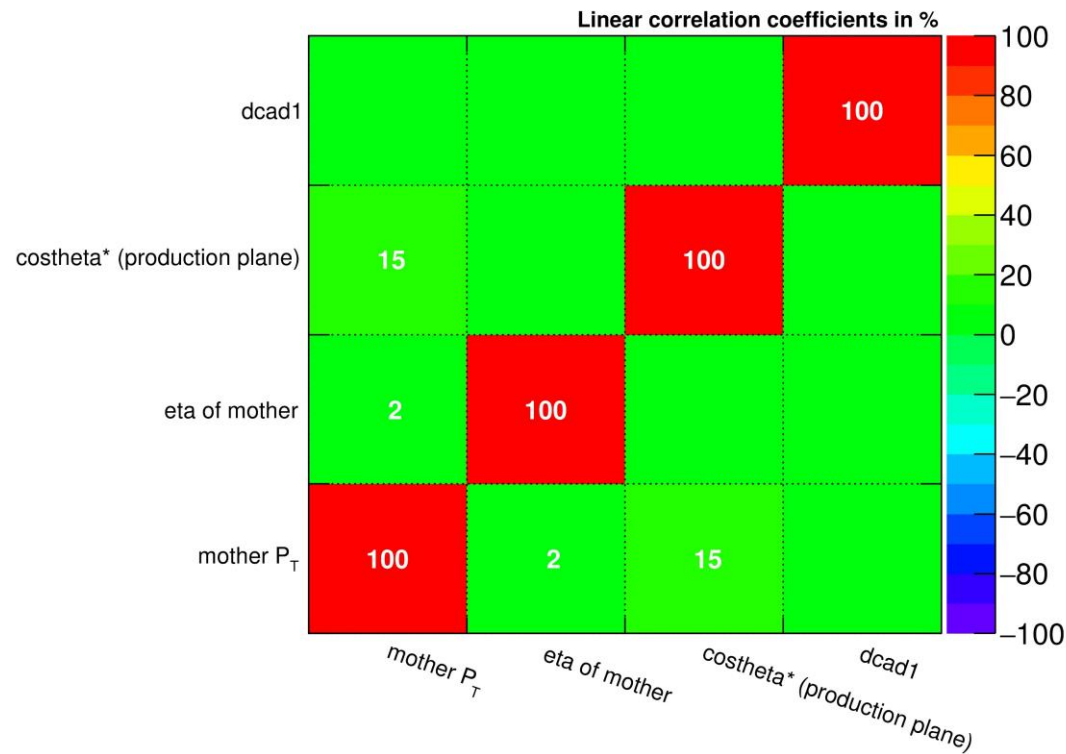
Correlation Matrix (background)



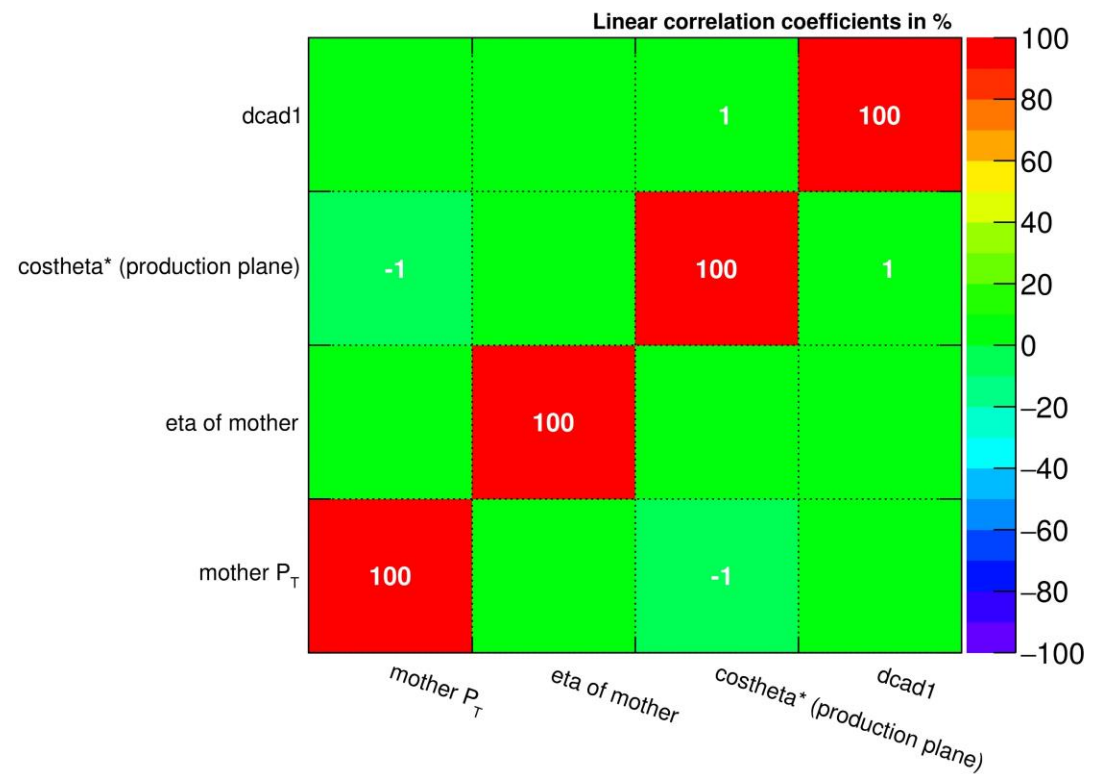
Input features used for training phase.

Selected features have less correlation with invariant mass

Correlation Matrix (signal)



Correlation Matrix (background)



DT algorithm

- An event is either classified as signal or background by either passing or not passing a condition (cut) on a specific node until a decision is made.
- In order to determine these conditions(cuts) , the decision tree is grown starting from the root node.

Where to stop ? Stop the splitting when we reach maximum tree depth or we have exhausted all our features!

Class of a leaf : If $\text{purity}(S/S+B) > 0.5$ then signal, otherwise background.

How do I select which feature I take first?

Gini Index : Define as $p^2 + q^2$ (at a given node)

- p : fraction of positive (signal) events
- $q=1-p$: fraction of negative (background) events

Weighted Gini Split:

$$I_G = f_{\text{left}} * G_{\text{left}} + f_{\text{right}} * G_{\text{right}}$$

where,

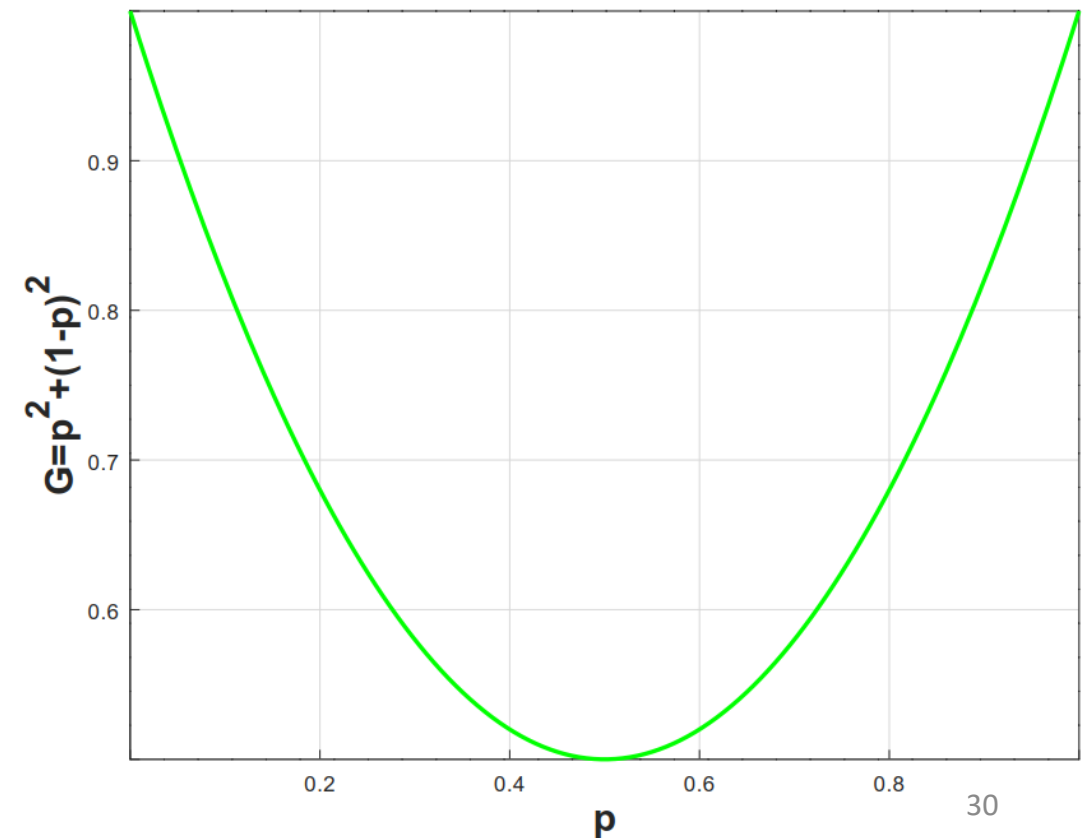
f_{left} = fraction of events which go in the left split

f_{right} = fraction of events which go in the right split

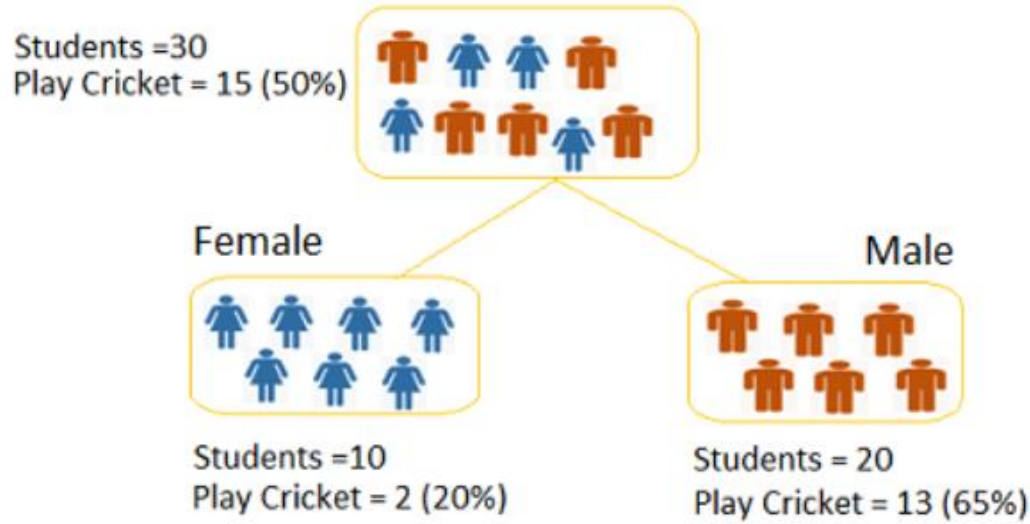
G_{left} = Gini index of left node

G_{right} = Gini index of right node

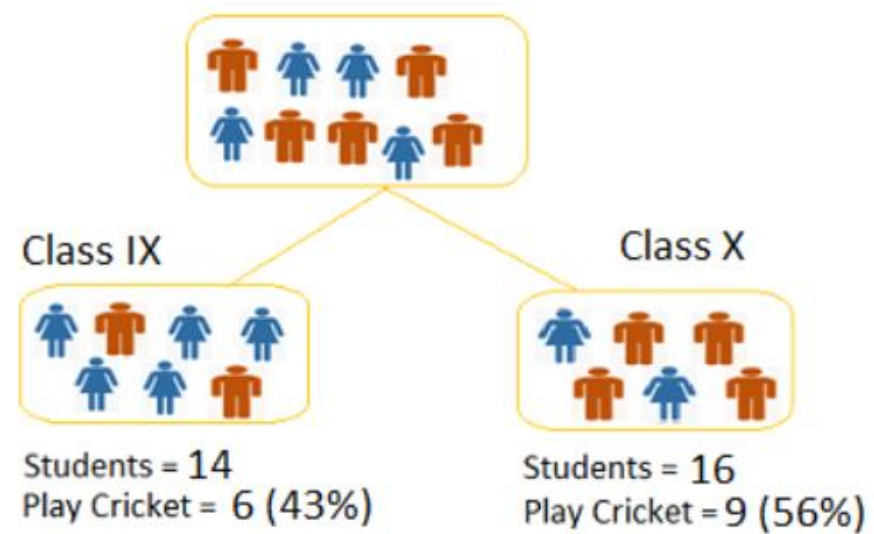
Plot of Gini index vs p (positive/signal)



Split on Gender



Split on Class



Split on Gender:

1. Gini for sub-node Female = $(0.2) * (0.2) + (0.8) * (0.8) = 0.68$
2. Gini for sub-node Male = $(0.65) * (0.65) + (0.35) * (0.35) = 0.55$
3. Weighted Gini for Split Gender = $(10/30) * 0.68 + (20/30) * 0.55 = 0.59$

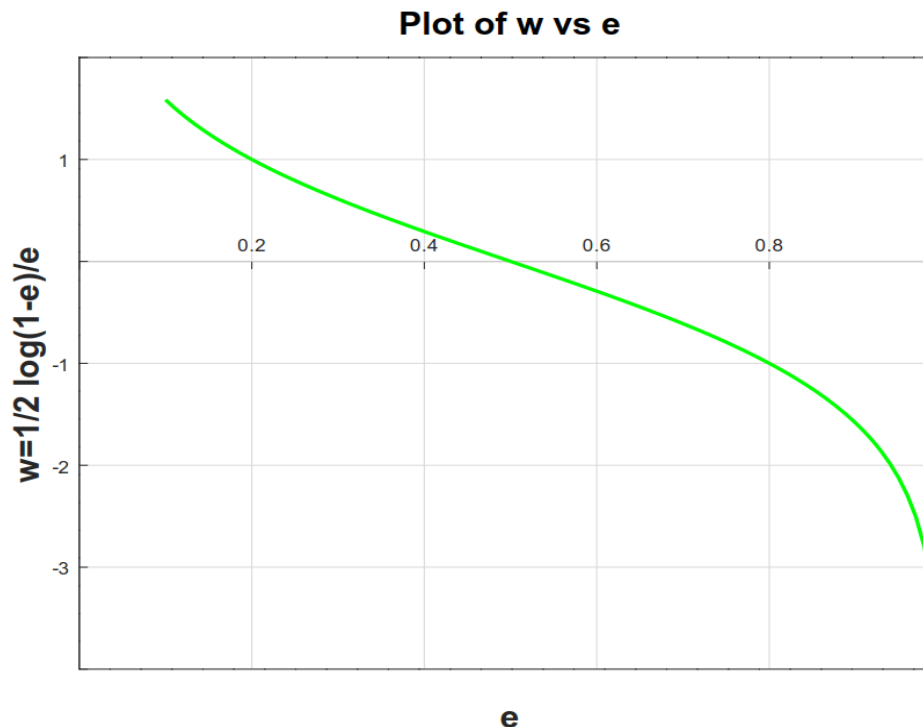
Similar for Split on Class:

1. Gini for sub-node Class IX = $(0.43) * (0.43) + (0.57) * (0.57) = 0.51$
2. Gini for sub-node Class X = $(0.56) * (0.56) + (0.44) * (0.44) = 0.51$
3. Weighted Gini for Split Class = $(14/30) * 0.51 + (16/30) * 0.51 = 0.51$



AdaBoost: Formula for computing coefficient \hat{w}_t of classifier $f_t(x)$

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$



Is $f_t(x)$ good?

weighted_error(f_t) on training data	$\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)}$	\hat{w}_t
0.01	$\frac{1 - 0.01}{0.01} = 99$	$\frac{1}{2} \ln(99) = 2.3$
0.5	$\frac{1 - 0.5}{0.5} = 1$	0
0.99	$\frac{1 - 0.99}{0.99} = 0.01$	$\frac{1}{2} \ln(0.01) = -2.3$

Yes

No

A terrible classifier!
But negative of that:
-f is awesome! 😊

AdaBoost: Formula for updating weights α_i

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) = y_i \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) \neq y_i \end{cases}$$

Did f_t get \mathbf{x}_i right?

Yes
No

$f_t(\mathbf{x}_i) = y_i$?	\hat{w}_t	Multiply α_i by	Implication
Correct	2.3	$e^{2.3} = 9.98$	Increases importance at x_i
Correct	0	$e^0 = 1$	Keeps importance same
Mistake	2.3	$e^{-2.3} = 0.1$	Decreases importance at x_i
Mistake	0	$e^{-0} = 1$	Keeps importance same

AdaBoost: learning ensemble

- Start same weight for all points: $\alpha_i = 1/N$

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$

- For $t = 1, \dots, T$

- Learn $f_t(\mathbf{x})$ with data weights α_i

- Compute coefficient \hat{w}_t

- Recompute weights α_i

- Normalize weights α_i

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) = y_i \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) \neq y_i \end{cases}$$

- Final model predicts by:

$$\hat{y} = \text{sign} \left(\sum_{t=1}^T \hat{w}_t f_t(\mathbf{x}) \right)$$

$$\alpha_i \leftarrow \frac{\alpha_i}{\sum_{j=1}^N \alpha_j}$$

Threshold split selection algorithm

- **Step 1:** Sort the values of a feature $h_j(\mathbf{x})$:
Let $\{v_1, v_2, v_3, \dots, v_N\}$ denote sorted values
- **Step 2:**
 - For $i = 1 \dots N-1$
 - Consider split $t_i = (v_i + v_{i+1}) / 2$
 - Compute classification error for threshold split $h_j(\mathbf{x}) \geq t_i$
 - Chose the t^* with the lowest classification error

Income