

Wids
Project
2025

Predicting Physics Observables at the LHC
using Machine Learning on Advanced
Collider data

By
Ankit Kumar
25N0292
M.Sc Physics

Abstract

In high-energy physics analyses involving heavy-flavor hadrons, separating prompt signals from non-prompt contributions and background noise signals is a critical task. Traditional cut-based approaches often fail to capture complex correlations among topological observables, motivating the use of multivariate techniques.

In this project, we study the classification of candidates into Prompt, Non-Prompt, and Background signal using Boosted Decision Trees (XGBoost). The analysis is performed in exclusive transverse momentum (p_T) intervals, ensuring p_T -dependent modeling and avoiding kinematic bias. Only topological variables are used for training, while kinematic variables are reserved for validation and post-training studies.

This report summarizes the dataset preparation, feature selection, model training strategy, hyperparameter tuning, and preliminary performance evaluation conducted up to the mid-term stage.

Introduction

Ultra-relativistic heavy-ion collisions at modern collider facilities such as the Large Hadron Collider (LHC) and the Relativistic Heavy-Ion Collider (RHIC) provide a unique opportunity to study strongly interacting matter under extreme conditions of temperature and energy density. These experiments aim to recreate, for fleeting moments, the quark–gluon plasma (QGP), a deconfined state of quarks and gluons believed to have existed in the early Universe. Since the QGP is extremely short-lived, its properties cannot be measured directly and must instead be inferred through a variety of experimental observables.

A major experimental challenge in heavy-ion and proton–proton collision analyses is the reliable separation of genuine signal particles from large combinatorial backgrounds, as well as the disentanglement of different production mechanisms. For example, inclusive J/ψ or D^0 yields receive contributions from both **prompt production**, originating directly at the primary interaction vertex, and **non-prompt production**, arising from the weak decays of longer-lived beauty hadrons. Traditionally, such separation has relied on template fitting techniques using invariant mass distributions and decay-length-based observables. While effective, these methods are often statistically limited, computationally expensive, and constrained in their ability to perform fine-grained, multidimensional analyses.

In recent years, **machine learning (ML) techniques have fundamentally transformed the way data are analyzed in heavy-ion and high-energy physics**. By learning complex, nonlinear correlations among multiple kinematic and topological observables, ML algorithms provide a powerful alternative to conventional cut-based or template-based approaches. ML-based approaches allow for faster, more flexible, and more robust extraction of physics observables,

while reducing reliance on model-dependent fitting procedures. As heavy-ion experiments move toward increasingly high-precision measurements and larger datasets, machine learning is expected to play an essential role in advancing our understanding of QCD matter under extreme conditions.

Motivated by these developments, the present project explores the application of machine-learning techniques to heavy-flavor classification, with particular emphasis on topology-based and kinematic observables using XGBoost. By building upon recent advances in ML-driven analyses, this work aims to demonstrate how modern data-driven methods can enhance signal identification, improve background rejection, and ultimately deepen our insight into heavy-ion collision dynamics.

Background

D^0 mesons are unstable particles created during high ion collisions and decays into:

$$D^0 \rightarrow K + \pi$$

D^0 itself is not seen. Only its decay products are seen. It is inferred by: finding two tracks (K and π), checking if they come from the same point (SV) or that point is displaced from PV and checking if their combined mass matches D^0 .

For every real D^0 there may be hundreds of fake combinations, two random tracks that *happen* to cross tracks from different particles, this is background noise.

D^0 that is detected at the primary vertex (the point where collision takes place) is termed as prompt signal. Apart from D^0 , a bottom quark is produced at the PV, which forms a B meson. The B meson decays as:

$$B \rightarrow D^0 + X$$

When this D^0 decays, a non-prompt signal is detected.

In this project, we will be using the following important topological variables (with column names in dataset) for the classification model training:

1. Cosine of pointing angle (fCpaD0) : Related to reconstructed momentum vector. The sum of momentum vec of daughters points towards PV
2. Transverse pointing (fCpaXYD0): Related to reconstructed momentum vector in XY plane
3. Distance from the PV (fDecayLengthXYD0)
4. Product of Daughter tracks' impact parameters (fImpactParameterProductD0)
5. Impact parameter of soft pion (from D^0) (fImpactSoftPi)
6. Maximum normalized impact parameter difference (fMaxNormalisedDeltaIPD0)

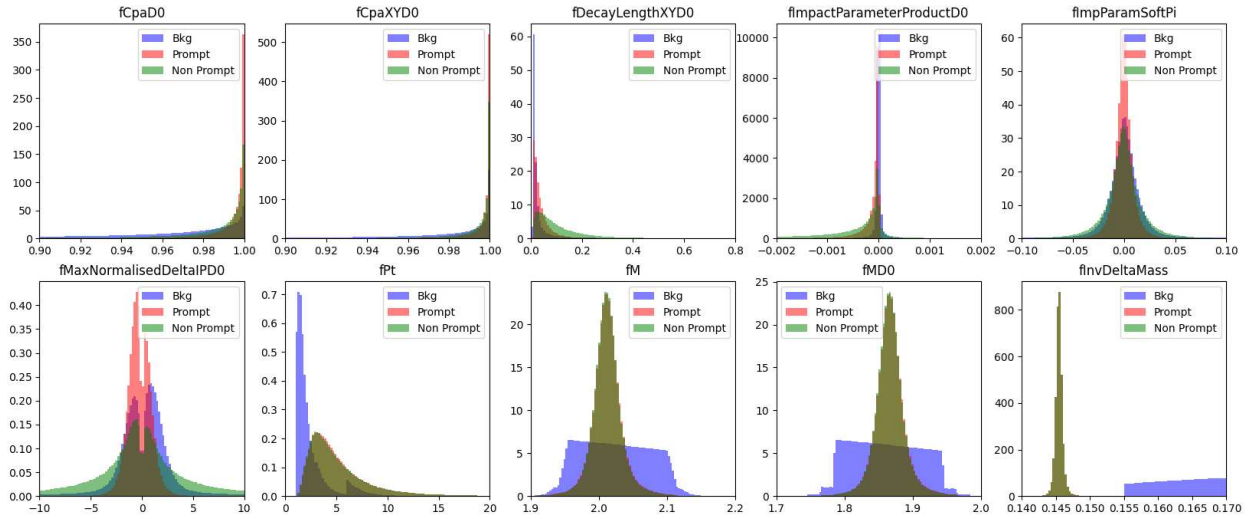
The following kinematic variables (with column names in dataset) useful for binning and verification of the three classes:

1. Transverse momentum of D^0 (fPt)
2. Invariant mass of reconstructed candidate (fM)

$$M = \sqrt{E^2 - |\vec{p}|^2}$$

3. Invariant mass with D^0 mass hypothesis (fMD0)
4. Invariant mass difference (D^*) (fInvDeltaMass)

$$\Delta M = M(K\pi\pi_{\text{soft}}) - M(K\pi)$$



Workflow

1. Dataset Description

Three ROOT files are used in this analysis:

- **Prompt_DstarToD0Pi.root**
- **Nonprompt_DstarToD0Pi.root**
- **Bkg_DstarToD0Pi.root**

Each file contains reconstructed D^* candidates stored in a common TTree structure. The datasets correspond to:

Class	Description
Prompt	D mesons produced directly in the primary interaction
Non-Prompt	D mesons from weak decays of B hadrons
Background	Combinatorial background candidates

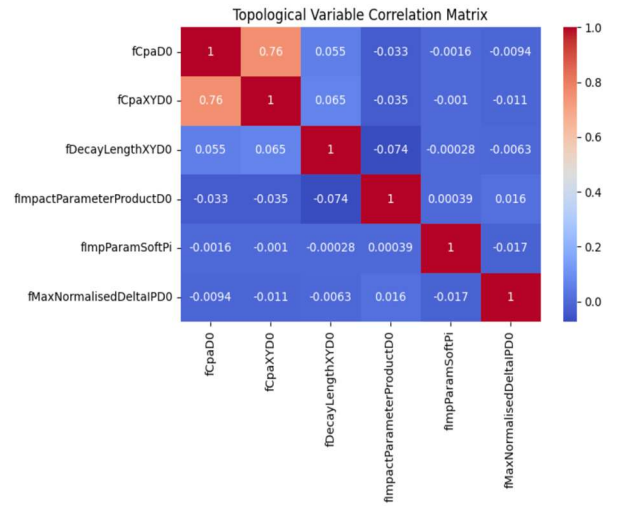
The ROOT files are accessed using the **uproot** Python library, enabling ROOT-independent data analysis. Further these are converted into Pandas data frames.

2. Feature Selection

3.1 Topological Variables (Used for Training)

Only topology-related observables are used as inputs to the BDT to ensure that the classifier learns genuine decay topology rather than kinematic differences.

These variables encode information about decay geometry, pointing angles, and impact parameters, which are known to be powerful discriminants between signal and background.



3.2 Kinematic Variables (Used Only for Plotting)

3. pT Binning Strategy

To study the p_T dependence of the classification performance and avoid kinematic bias, the analysis is performed in exclusive p_T intervals defined by: [1, 1.5, 2, 3, 4, 6, 8, 10, 12, 16, 24] GeV/c. Here, we have chosen [4,6] GeV/c bin.

4. Train–Test Split

The binned dataset is first divided into **training** and **testing** samples using a **70% / 30% split**:

- Training fraction: 0.7
- Testing fraction: 0.3

The number of instances for each class must be equal in the training set in order to avoid any bias towards the majority class. To ensure this, the trained data is oversampled for whichever class with lesser instances. We have avoided under sampling so as not to lose data.

5. Model Choice

Boosted Decision Trees (BDT)

Boosted Decision Trees are chosen due to their:

- Ability to capture non-linear correlations
- Robust performance with limited feature sets
- Long-standing validation in heavy-flavor analyses

The **XGBoost** implementation is used, configured for **multi-class classification** with three output classes (Background, Prompt, Non-Prompt).

6. Hyperparameter Tuning Strategy

Random search method is used for the tuning of the hyperparameters.

Parameter Grid:

```
"max_depth":    [2, 4, 6],  
"learning_rate": [0.3, 0.5, 0.7],  
"n_estimators":  [400, 600, 1000],  
"subsample":    [0.7, 0.8, 1],  
"colsample_bynode": [0.5, 0.7, 1.0],  
"min_child_weight": [6, 7, 8],  
"colsample_bytree": [0.5, 0.7, 1.0],  
"reg_alpha":     [0, 0.001, 0.01],  
"reg_lambda":    [0, 0.001, 0.01],  
"gamma":         [0, 0.1, 0.5],
```

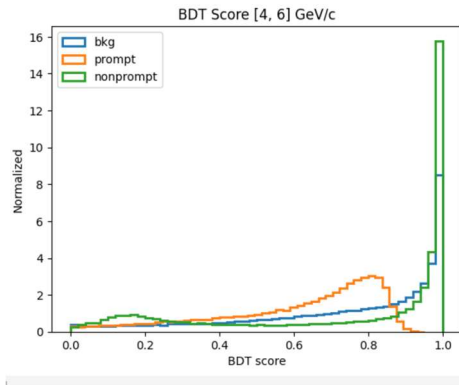
With 5 combinations from the above grid and 8 cross validations, the model with optimum ROC AUC number is chosen. The obtained model has the following hyperparameter

Best Parameters: {'subsample': 0.7, 'reg_lambda': 0, 'reg_alpha': 0, 'n_estimators': 400, 'min_child_weight': 6, 'max_depth': 6, 'learning_rate': 0.3, 'gamma': 0.1, 'colsample_bytree': 1.0, 'colsample_bynode': 0.5}

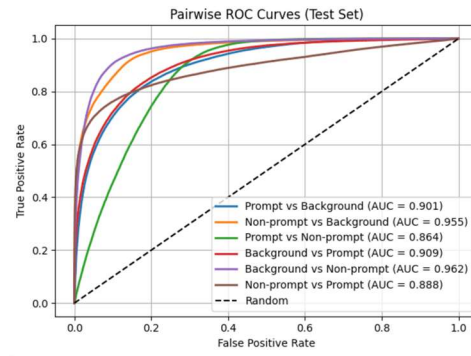
Best Accuracy: 0.920464188024964

7. Performance Evaluation

7.1 BDT Score Distributions



7.2 ROC Curves



Conclusion

This mid-term work establishes a robust, pT-differential machine learning framework for heavy-flavor signal classification using topology-based observables. The methodology follows best practices in high-energy physics and provides a strong foundation for further extensions involving PID information and real data applications.

References

1. [A short course on Relativistic Heavy Ion Collisions](#) by [A. K. Chaudhuri](#)
2. [Relativistic Kinematics](#) by [Raghu Nath Sahoo](#)
3. [Inclusive, prompt and nonprompt \$J/\psi\$ identification in proton-proton collisions at the Large Hadron Collider using machine learning](#) by Suraj Prasad ,[†] Neelkamal Mallick ,[‡] and Raghu Nath Sahoo