

Selection of $B_s^0 \rightarrow \psi(2S)K_S^0$ decays via multivariate analysis

Lukas Bertsch

lukas.bertsch@tu-dortmund.de

Tabea Hacheney

tabea.hacheney@tu-dortmund.de

Tom Troska

tom.troska@tu-dortmund.de

Start of course: 13th of June 2024

TU Dortmund University – Faculty of Physics

Contents

1. Motivation	3
2. Theory	3
3. The LHCb detector	4
4. Analysis strategy	6
5. Analysis	6
5.1. Definition of a signal window	7
5.2. Feature selection	7
5.3. Training of a multivariate classifier	9
5.4. Optimization of the classification threshold	10
5.5. Evaluation of the signal yield	11
6. Discussion	13
References	14
A. Appendix	15
A.1. Correlations and distributions of the variables used for the MVA	15

1. Motivation

During Run 2 (2015-2018) of the LHCb experiment at CERN, proton bunches with a centre of mass energy of 13 TeV collided approximately 40×10^6 times per second [1]. Next to potentially interesting processes, a lot of very common decays as well as background is recorded. To identify and extract rare decay channels, an extensive analysis including machine learning methods is needed. This analysis is aimed at finding events coming from the decay $B_s^0 \rightarrow \psi(2S)K_S^0$ using data recorded by the LHCb experiment in Run 2.

2. Theory

Before diving into the analysis, one has to understand the underlying kinematics of the decay. The B_s^0 meson consists of a s and a \bar{b} quark (next to other sea quarks and gluons) and hence has a neutral net electric charge. As for this decay channel, the b quark of the B_s^0 meson interacts with a W^+ gauge boson and results in a charm pair $c\bar{c}$ and a \bar{d} quark in the final state. Together with the strange quark from the B_s^0 meson, the final state consists of a $\psi(2S) [c\bar{c}]$ meson and a K_S^0 ("k-short") being a superposition of $d\bar{s}$ and $\bar{d}s$. This decay can only occur via the weak force, since the quark flavour changes from \bar{b} to \bar{d} . The Feynman diagram depicting this process is shown in Figure 1.

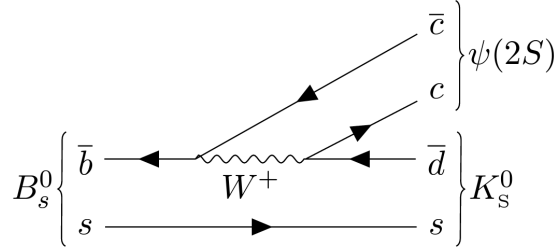


Figure 1: Feynman diagram of $B_s^0 \rightarrow \psi(2S)K_S^0$ [1].

The $\psi(2S)$ and K_S^0 are not directly detected in the detector. While $\psi(2S)$ mainly decays via strong interactions, it also decays via electromagnetic force into a lepton pair (l^+l^-). Since the detector has a good efficiency for detecting muons, only muon pairs will be regarded ($\mu^+\mu^-$). Due to the conservation of energy and momentum, the total mass of the muon pair is approximately equal to the mass of the $\psi(2S)$ meson. With cuts on the muon pair masses, the collected data can already be reduced to lesser signal candidates while explicitly rejecting other charm resonances $c\bar{c}$.

The dominant decay channel of K_S^0 is a pair of pions ($\pi^+\pi^-$) via weak interaction. The K_L^0 meson also consists of the superposition of $d\bar{s}$ and $\bar{d}s$ and has the same decay channels. However, the branching fraction of that decay channel is two magnitudes smaller than for the K_S^0 and will therefore not interfere with the relevant signal candidates to a significant level. The mass of the K_S^0 can be reconstructed with the produced pion pairs.

The given recorded data consists of various variables including kinematic properties of

the muons, pions as well as the reconstructed $\psi(2S)$ and K_S^0 mesons. The B_s^0 can be reconstructed using the properties of the $\psi(2S)$ and K_S^0 while using their true invariant masses, being $m_{\psi(2S)} = (3686.097 \pm 0.011) \text{ MeV}$ and $m_{K_S^0} = (497.611 \pm 0.013) \text{ MeV}$ [2]. The recorded data set does not only include $B_s^0 \rightarrow \psi(2S)K_S^0$ signal. Moreover, it primarily contains combinatorial background, consisting of random particle tracks whose particle trajectories align with the signal decay. The much more abundant decay $B^0 \rightarrow \psi(2S)K_S^0$ is also included in the recorded data. To differentiate between background and signal, an extensive analysis including training a classifier is needed.

To evaluate and optimize the classifier (maximizing the number of signal candidates while keeping the background low), a loss function is to be minimized. The inverse of such loss functions is called figure of merit (FOM). In this analysis, the Punzi figure of merit for optimizing the signal efficiency

$$\text{FOM} = \frac{\epsilon_{\text{sig}}}{5/2 + \sqrt{N_{\text{bkg}}}} \quad (1)$$

is to be maximized with N_{bkg} being the background candidates in the signal region (combinatorial background + $B^0 \rightarrow \psi(2S)K_S^0$ events) and ϵ_{sig} being the efficiency of classifying the signal. The signal efficiency is calculated by using a Monte Carlo simulation of $B_s^0 \rightarrow \psi(2S)K_S^0$. The significance in the signal region can be calculated by

$$m = \frac{N_{\text{sig}}}{\sqrt{N_{\text{sig}} + N_{\text{bkg}}}}, \quad (2)$$

using the given number of signal and background candidates by the classifier.

To train an efficient and well performing classifier, the right choice of variables is needed. One has to look for those variables, which differ the most for background and signal. To find these, one can calculate the largest distance between the cumulative probability distributions F^i of these variables:

$$\sup_n |F_n^1 - F_n^2|, \quad (3)$$

with the index n running over all bins of the distributions.

3. The LHCb detector

Located at CERN, the LHC is the largest particle collider in the world and houses a number of experiments. Opposing protons beams are collided at four different interaction points. At one of these, the LHCb experiment is installed. This detector is constructed as a single-arm forward spectrometer which allows for an optimized detection in the pseudorapidity range of $2 \leq \eta \leq 5$. With the help of this design, the LHCb detector is well suited for the detection of decays that include b and c quarks because the production of these is favored for small angles along the beam axis.

The LHCb experiment consists of several sub-detectors, each with a different special purpose. A cross-section of the detector used during the second data-taking run is depicted in Figure 2. In the immediate proximity of the interaction point, the Vertex

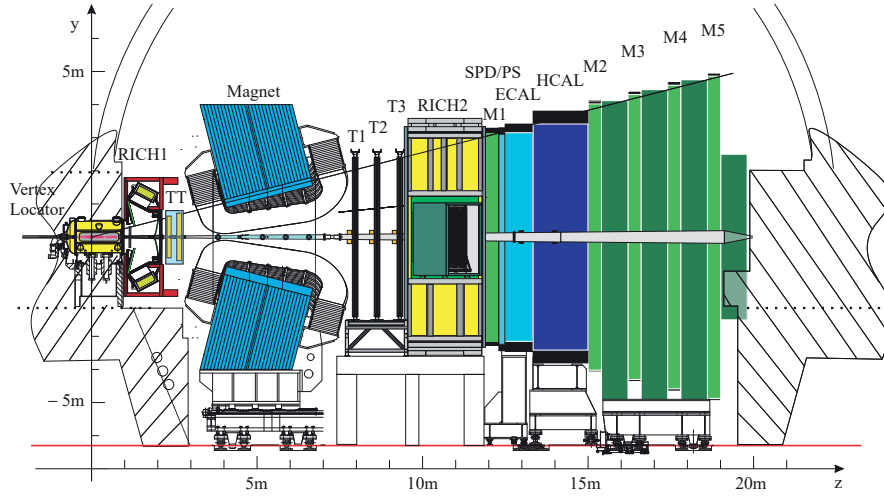


Figure 2: Cross-section of the LHCb detector. A right-handed coordinate system with the z-axis parallel to the beam pipe and the y-axis orientated to the top is used to describe the detector. The different sub-detectors and their purpose and functionality are discussed in this section [3].

Locator (VELO) is located. Primary and secondary vertices are identified by the 52 silicon pixel modules. The detector is divided into two halves that are retractable. During the injection phase, the VELO is sitting in its retracted position because the beam is not as focused as during the stable beam phase. Due to the closeness of the VELO to the interaction point, the primary vertices can be reconstructed with great accuracy.

One of the two Ring Imaging Cherenkov detectors (RICH) is located right before the 1.4 T dipole magnet while the other RICH detector is situated behind the tracking stations. Two materials with different refractive indices are used in the two RICH detectors to determine the velocity of the particles via the Cherenkov effect. The implementation of two detectors allows for a larger momentum range reconstruction.

The position of the particles are detected by the tracking stations, the Tracker Turicensis (TT) and T1-T2. With this information and the previously measured velocity, the momentum of the particles can be calculated.

For the determination of the energy of the particles, the electromagnetic and hadronic calorimeters (ECAL and HCAL) are employed. To accomplish this task for the ECAL, the shashlik calorimeter technology is used. Here, an absorber material and a detector layer of a scintillating material are stacked alternately. The HCAL is constructed similarly but with the scintillating tiles running parallel to the beam axis.

Muon stations at the end of the detector identify the muons that are passing through the detector. This is accomplished by four multi-wire proportional chambers.

Because of the large number of potential events, a trigger system is used to identify

interesting decays. The data of these are then saved for further offline analysis.

4. Analysis strategy

The data used in this analysis consists of three different data samples. One is the actual measured data, which contains reconstructed $B^0 \rightarrow \psi(2S)K_S^0$ candidates and is dominated by combinatorial background and a B_d^0 peak. The two other datasets are Monte Carlo simulation samples of the signal decay $B_s^0 \rightarrow \psi(2S)K_S^0$ and the control channel decay $B_d^0 \rightarrow \psi(2S)K_S^0$.

First, the training samples for the background and the signal need to be defined. In the case for the data of $B_d^0 \rightarrow \psi(2S)K_S^0$ decays, the background is predominately combinatorial. The recorded LHCb data is used for the training of this classifier. It is to note that the region, where most of the signal is expected, needs to be excluded from this to not get any signal bias into the background classification.

The training samples for the signal classification cannot be sourced from the recorded data but must be simulated with Monte Carlo simulations. This leads to the problem of imperfect simulation data that does not match the true data of the decay. This is due to not perfect theoretical models and computing constraints. The challenges of this problem can at least be partially overcome by introducing computing weights. The here provided simulation data contain a variable called `kinematic_weights` that helps to correct for the mismatches between simulation and recorded data.

Feature selection is another important aspect of this analysis. The classifier uses a number of features to decide whether the data is background or signal. Features, that can be used here, are extracted from the control channel $B^0 \rightarrow \psi(2S)K_S^0$. This decay is similar to the decay $B_s^0 \rightarrow \psi(2S)K_S^0$ and therefore its features can be fed to the classifier. In principle, any number of features can be used by a classifier but the runtime scales with the number of features and too many features will eventually lead to overfitting. Some of the here provided features are also redundant and do not include additional information. For the purpose of ascertaining the best classification threshold, the FOM (Equation 1) is determined by calculating the signal efficiency and number of background events in the signal region.

In the final step, the trained and optimized classifier is used to classify the recorded LHCb data. This leads to the removal of the combinatorial background so that the B^0 and B_s^0 peaks are clearly recognizable. The number of signal events can then be extracted.

5. Analysis

This analysis uses measured data with $B^0 \rightarrow \psi(2S)K_S^0$ candidates and Monte Carlo simulations of $B_s^0 \rightarrow \psi(2S)K_S^0$ and $B_d^0 \rightarrow \psi(2S)K_S^0$ decays. In Figure 3, the invariant mass distribution of the reconstructed B_s^0 events in the signal simulation is shown. As expected, a clear peak at the B_s^0 mass can be seen. The mass distribution of the reconstructed B^0 candidates in real data can be seen in Figure 4. Here, a peak at the nominal B_d^0 mass can be seen. However, due to dominating combinatorial background,

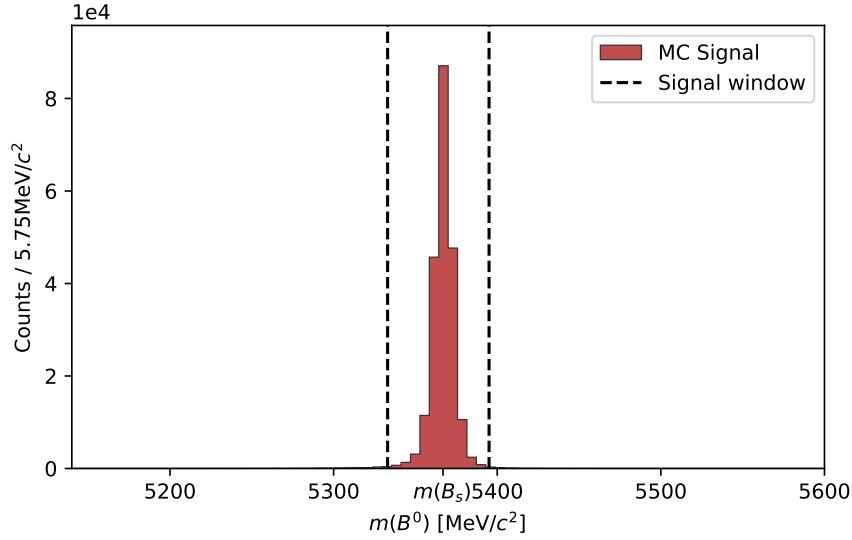


Figure 3: Invariant mass distribution of the B_s^0 candidates for the signal channel simulation data.

no peak at the B_s^0 mass is visible. The dataset also contains $sWeights$ which can be used to extract the contribution of the control channel in the dataset. By weighting the mass histogram with the $sWeights$, only the control channel contribution remains in the plot, as can be seen in Figure 5.

5.1. Definition of a signal window

The mass distribution of the signal decay peaks only in a short window of the whole mass range given in the dataset. In order to know, where signal is expected, a signal window has to be defined. This is done by calculating every interval containing 99% of data in the signal simulation and choosing the shortest interval. Here, this interval follows as 5333.4 to 5394.6 MeV/c^2 , defining the signal window of 5333 to 5395 MeV/c^2 which can also be seen in the aforementioned plots (3, 4). Subsequently, the upper sideband (‘ USB ’), containing mostly combinatorial background, is defined as the area with reconstructed mass $> 5400 \text{ MeV}/c^2$.

5.2. Feature selection

In order to train a multivariate classifier capable of separating signal from background, meaningful features from all available variables in the dataset have to be extracted. In total, 863 variables are listed in the dataset. After removing event, utility, trigger and spatial coordinate variables, 398 variables remain. For these variables, the correlation to the invariant mass is calculated and variables having a correlation coefficient of 0.3 or higher are excluded. Variables with too high correlation to the invariant mass would

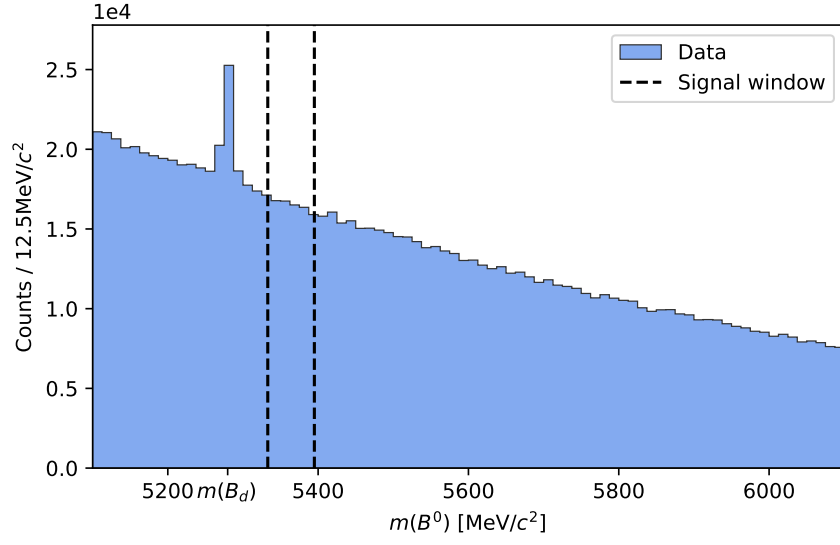


Figure 4: Invariant mass distribution of the B^0 candidates for the recorded LHCb data.

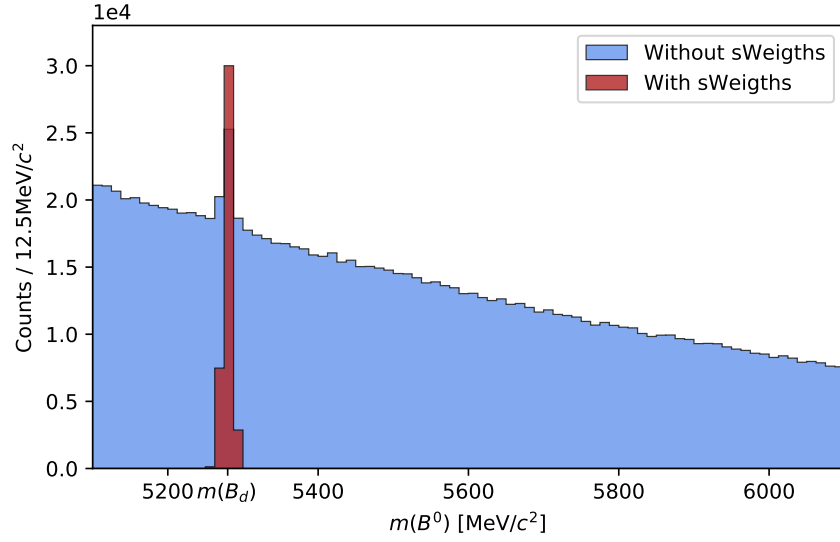


Figure 5: Invariant mass distribution of the B^0 candidates for the recorded LHCb data with and without the sWeights.

introduce a bias in training the classifier and could not be used in similarity checks between simulation and data, because the sWeights are based on the B^0 candidate mass. The remaining 390 variables are checked to be correctly modelled by simulation and have significantly different distributions for signal and background. This is done using the Kolmogorov Smirnov test statistic defined in Equation 3 for weighted distributions as a measure of similarity. To check agreement between simulation and data, the distributions of sWeighted data are compared to the control channel distributions and 190 variables with a test statistic of $d > 0.05$ are removed. The variables are also required to have a test statistic of $d > 0.2$ when comparing signal (simulation) and background (USB), leaving 94 variables. After further removing variables that are highly correlated (correlation > 0.9) to others or duplicates of another variable, 18 variables are left that are used for the training of the multivariate classifier. The distributions of four of these variables

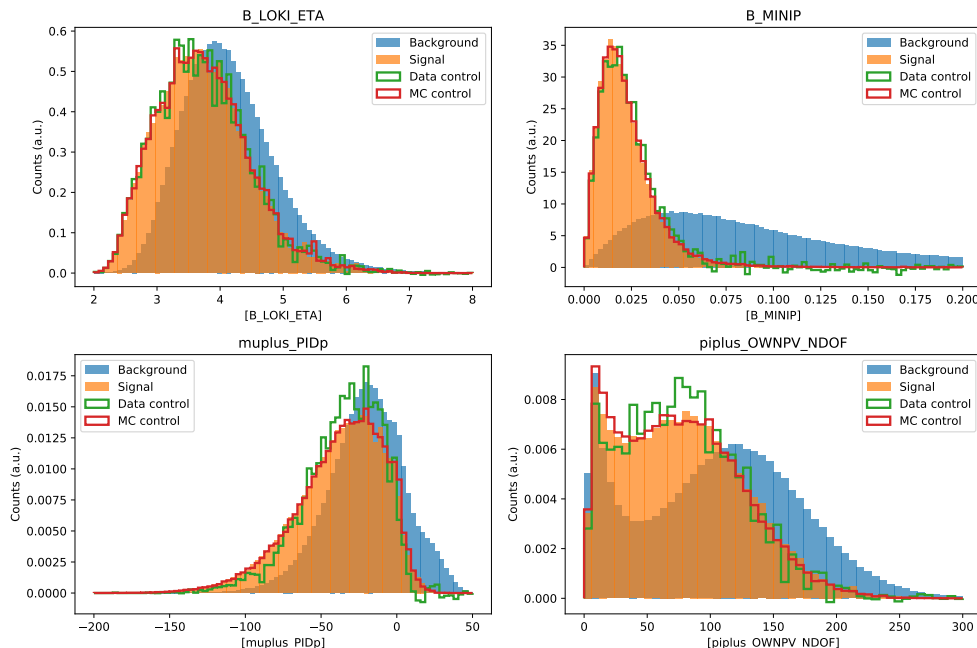


Figure 6: Distributions of four selected variables used in the MVA for simulation, reweighted data and background.

for signal and background, as well as control channel data and simulation are shown in Figure 6. Distributions for all variables and the correlation matrix can be found in the appendix A.1.

5.3. Training of a multivariate classifier

With the now selected variables, a multivariate classifier can be trained. For this purpose, a boosted decision tree as implemented in the package XGBoost [4] is used. The training data consists of 637 410 background events and kinematically reweighted signal simulation,

corresponding to 155 805 events. Because of this imbalance, each background event is weighted with a factor of 155805/637410. The hyperparameters of the classifier are optimized via a random search followed by a grid search using cross validation and a subset of 40 000 training samples. The resulting parameter values can be read from Table 1. For the classification, five individual BDT's are trained using 5-fold cross

Table 1: Hyperparameter values of the trained classifiers determined by grid search.

Parameter	Value
n_estimators	1000
learning_rate	0.1
max_depth	4
reg_lambda	1
n_iter_no_change	5

validation and the hyperparameters from Table 1. To evaluate the classifiers performance, the ROC curve is viewed and the area under the curve, as well as the accuracy are calculated. To check for overtraining, the response of the classifier is compared for training and simulation data in a logarithmic plot. The ROC curve and the train-test comparison of the fifth trained classifier can be seen in Figure 7. Additionally, the feature

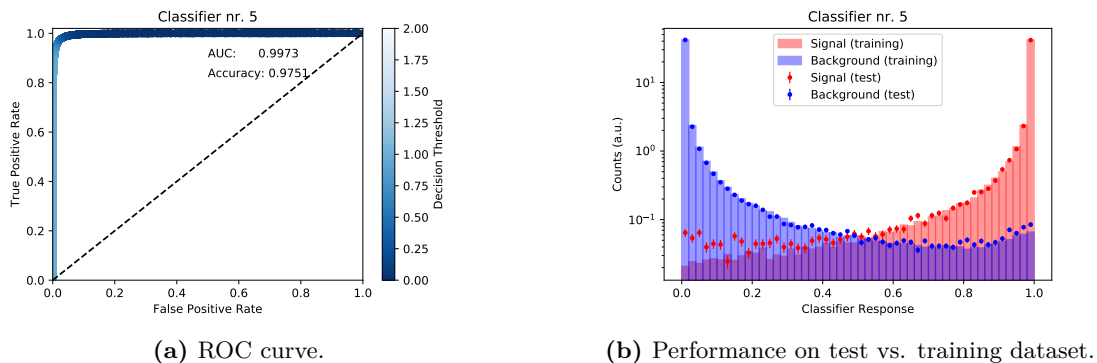


Figure 7: ROC curve (left) and response on training and test dataset (right) for one of the trained classifiers.

importance of the BDT variables is checked for imbalances. As can be seen in Figure 8, all variables are of similar importance.

5.4. Optimization of the classification threshold

After applying all BDT's to the data, a cut on the classifiers response has to be made. Therefore the mean classifier response between all 5 BDT's is calculated. The classification threshold separating signal and background is optimized using the Punzi figure of merit Equation 1. The signal efficiency ε is therefore calculated for each threshold as the selection efficiency on the signal simulation. The background yield in the signal region B

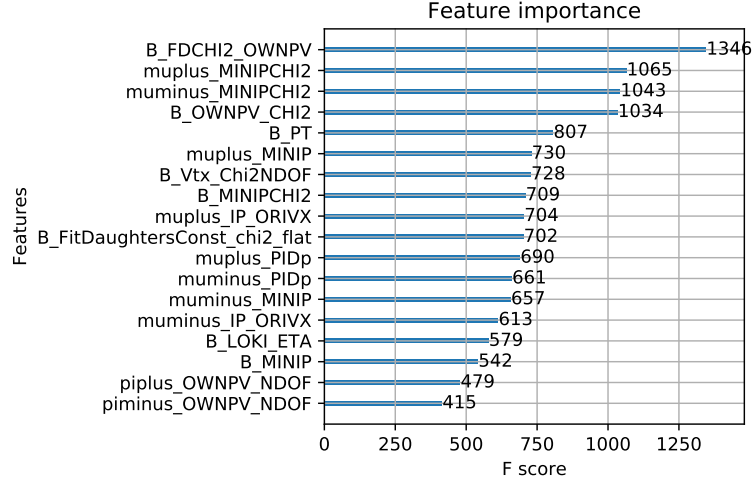


Figure 8: Feature importance of the BDT training variables.

is estimated via a fit of an exponential function to the upper sideband of the data after applying the selection threshold, and extrapolated to the signal window. For the fitting, the python library `iminuit` [5] is used. The resulting values of the figure of merit are plotted against different thresholds in Figure 9. The optimal classification threshold is

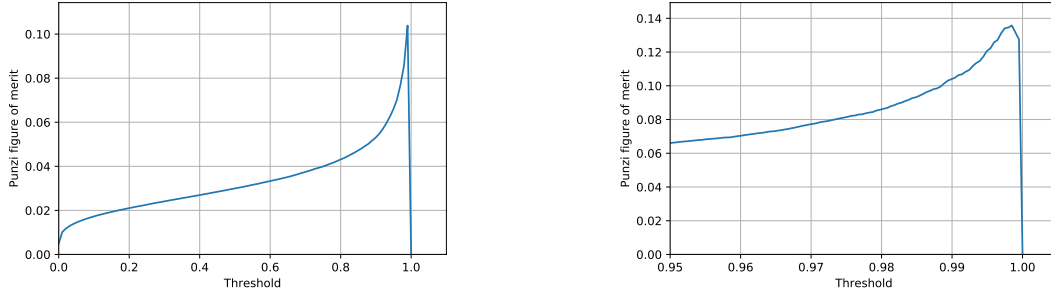


Figure 9: The Punzi figure of merit for the mean classifier response in different intervals of the threshold.

the maximum value of the Punzi figure of merit and reads as $t = 0.998$.

5.5. Evaluation of the signal yield

The invariant B^0 mass distribution of the data after applying the cut on the mean BDT response can be seen in a semi-logarithmic plot in Figure 10. A peak of the signal decay $B_s^0 \rightarrow \psi(2S)K_S^0$ can be seen.

In order to determine the signal yield, a fit to the mass spectrum of the B^0 candidates is applied. The signal peaks are modeled with two gaussian distributions, which start

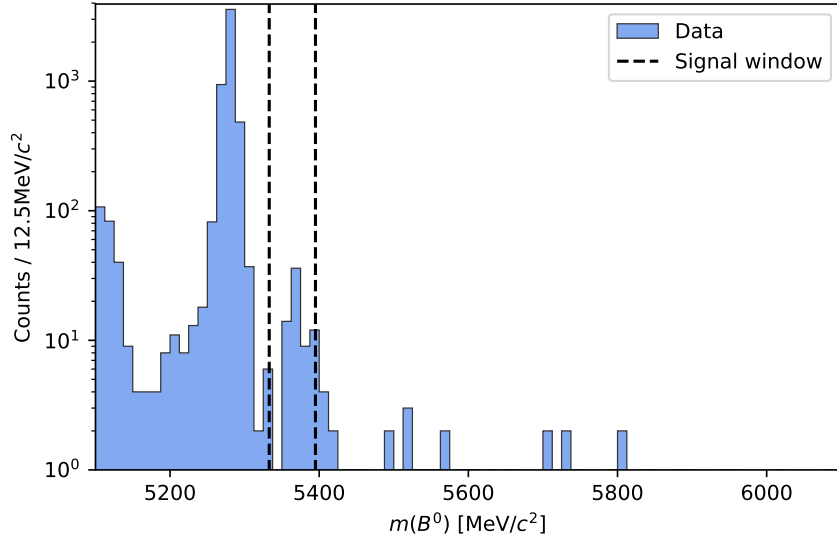


Figure 10: Semi logarithmic invariant mass distribution of the B^0 candidates in data, after the cut on the classifier response is applied.

values for the mean μ and width σ are defined by fits to simulation. The values itself are allowed to vary freely to account for mass resolution differences in data and simulation. The background is again modeled by an exponential function with decay constant τ . An extended, unbinned negative Log-Likelihood fit is performed, including the fractions s_{B_d} , s_{B_s} and b of signal, control channel and background components. The resulting graph is shown in Figure 11. The fit parameters follow as

$$\begin{aligned}
 s_{B_s} &= 43 \pm 8 & s_{B_d} &= 5010 \pm 71 \\
 \mu_{B_s} &= (5366.9 \pm 0.9) \frac{\text{MeV}}{c^2} & \mu_{B_d} &= (5279.9 \pm 0.9) \frac{\text{MeV}}{c^2} \\
 \sigma_{B_s} &= (4.6 \pm 0.8) \frac{\text{MeV}}{c^2} & \sigma_{B_s} &= (6.22 \pm 0.07) \frac{\text{MeV}}{c^2} \\
 b &= 492 \pm 24 & \tau &= 124 \pm 6.
 \end{aligned}$$

From these, the significance proxy as defined in Equation 2 of the observation of the signal can be calculated. Therefore, the signal and background events in the signal window are interpolated from the fit results as $n_{\text{sig}} = 43$ and $n_{\text{bkg}} = 30$. The significance proxy then reads $m = 5.07$.

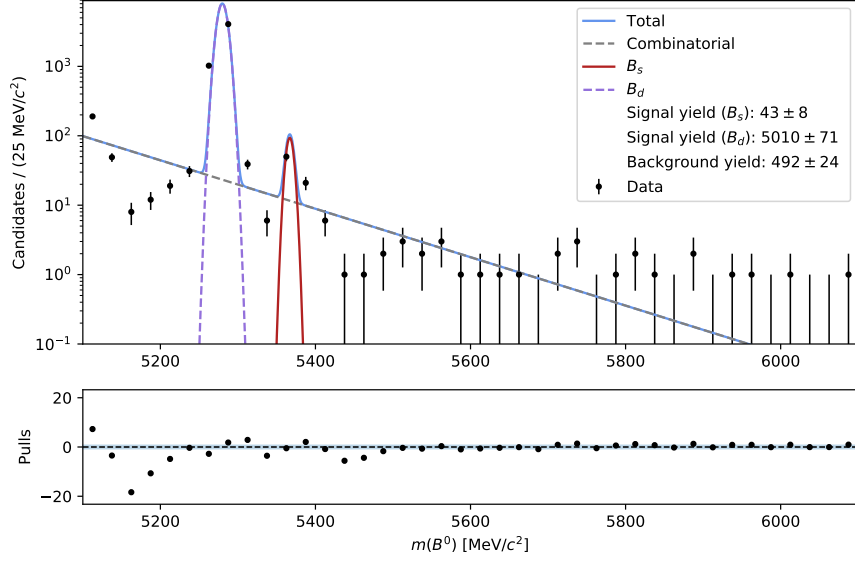


Figure 11: Fit to the invariant mass spectrum of the data in semi logarithmic depiction.

6. Discussion

The analysis of the data used for the study of the decay $B_s^0 \rightarrow \psi(2S)K_S^0$ yields a number of signal events of $n_{\text{sig}} = 43$ and a number of background events of $n_{\text{bkg}} = 30$ with a significance proxy of $m = 5.07$. Several aspects need to be considered for the validity of these numbers.

First, the BDT seems to be slightly overfitted as can be seen in Figure 7b. Here, marginally more signal events are classified in the test data as in the training data. This could be overcome with further training of the BDT by implementing and optimizing different hyperparameters.

Additionally, the signal of the B_s^0 decay is substantially smaller than the signal of the B_d^0 decay and in the lower energy spectrum less background could be removed. A possible explanation for this is that only the upper sideband was used for the training of the classifier. By adding data of the lower sideband or using variables that are less correlated to the mass, an enhanced background reduction could potentially be achieved.

A further aspect, that could be improved, is the fit model of the signal peaks. In general, these peaks are not symmetric but in this analysis a double gaussian fit is performed. A more complex model has the potential to optimize the fit.

For the significance proxy of about five sigma, it is to note that uncertainties were not considered for the calculation of this value so that the proxy likely overestimates the significance.

References

- [1] *Selection of $B0s \rightarrow (2S)K0 S$ events*. TU Dortmund, Faculty of Physics, AG Albrecht. 2020. URL: <https://moodle.tu-dortmund.de/mod/resource/view.php?id=1598013>.
- [2] *The Review of Particle Physics (2024)*. Particle Data Group. URL: <https://pdglive.lbl.gov/Viewer.action> (visited on 12/06/2024).
- [3] A. Augusto Alves Jr. et al. ‘The LHCb Detector at the LHC’. In: *JINST* 3 (2008), S08005. DOI: 10.1088/1748-0221/3/08/S08005.
- [4] T. Chen and C. Guestrin. ‘XGBoost: A Scalable Tree Boosting System’. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD ’16. San Francisco, California, USA: ACM, 2016, pp. 785–794. DOI: 10.1145/2939672.2939785.
- [5] H. Dembinski and P. Ongmongkolkul et al. ‘scikit-hep/iminuit’. In: (Dec. 2020). DOI: 10.5281/zenodo.3949207.

A. Appendix

A.1. Correlations and distributions of the variables used for the MVA

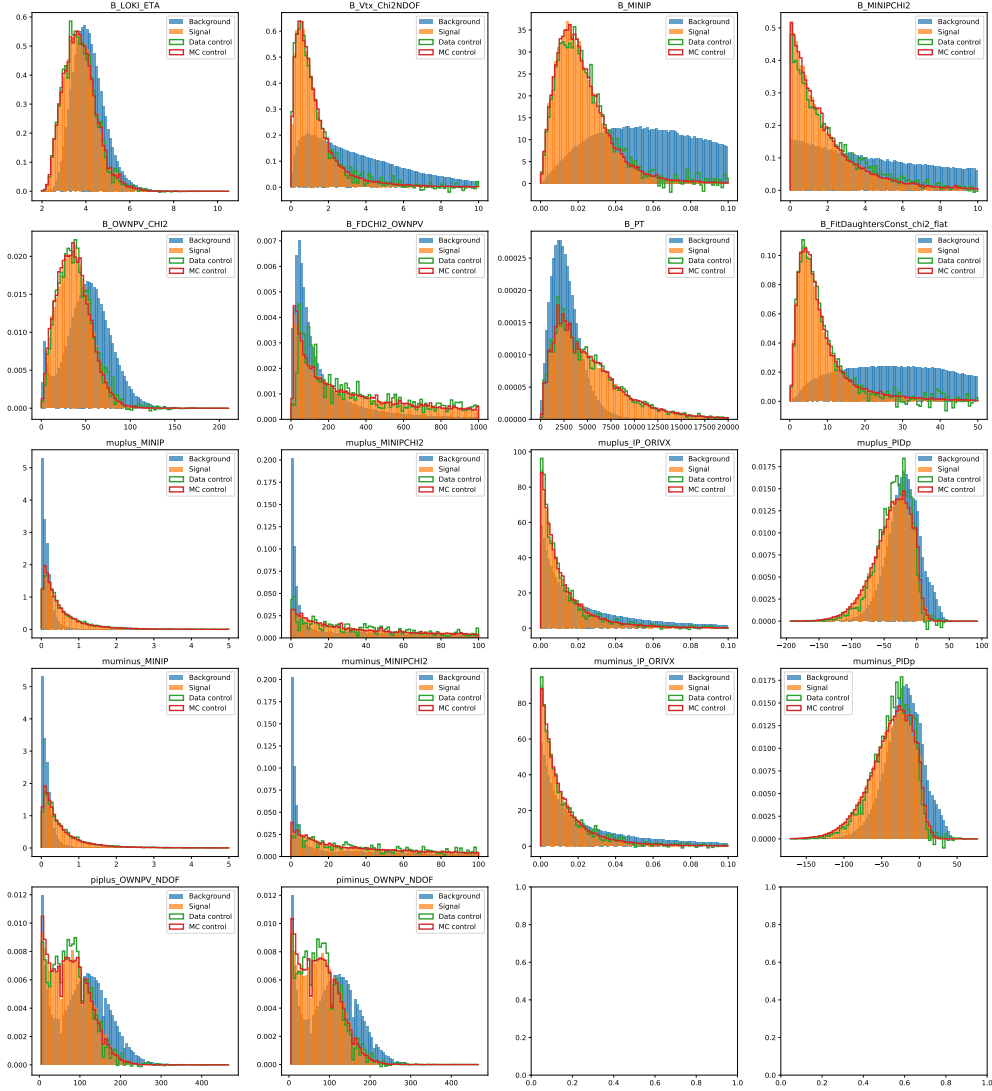


Figure 12: Distributions of the variables used in the MVA for simulation, reweighted data and background.

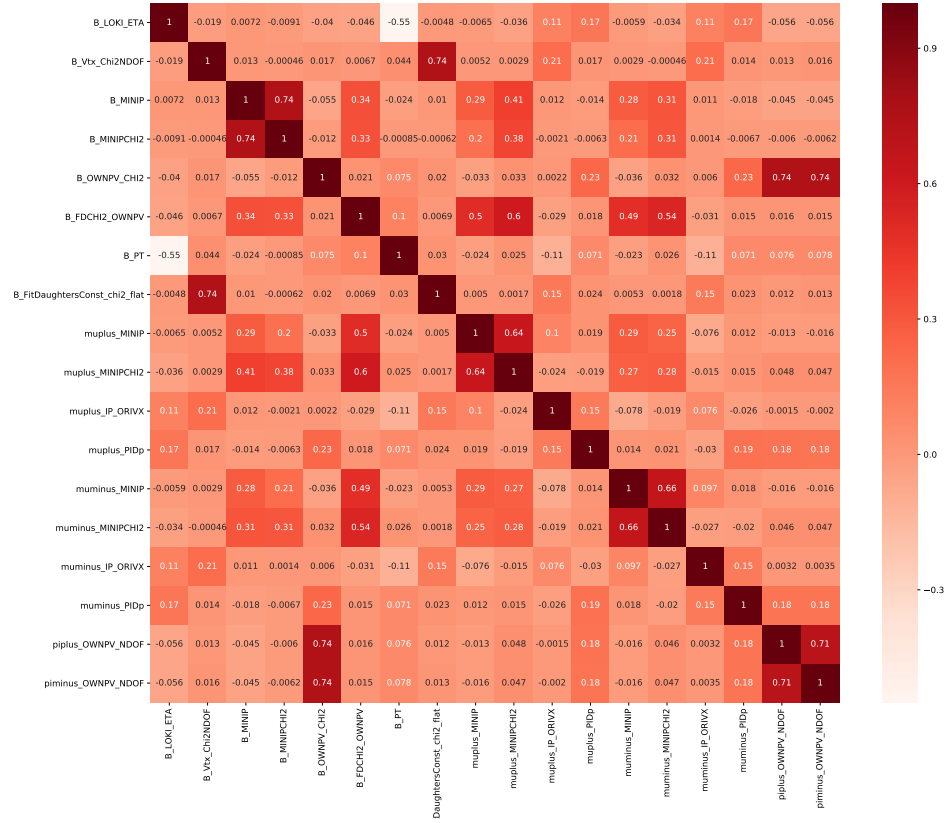


Figure 13: Correlations between the selected variables.