A multivariate selection of charmless B^+ decays to three mesons in LHCb detector

> J. Helder Lopes IF-UFRJ

for the UFRJ-CBPF LHCb BnoC Working Group

July, 31th 2018

1/16

Overview

- Motivation for the analysis
- Motivation for the MVA
- MVA in Run I
- MVA in Run II
 - Update from Run I
 - New reduced set of variables
 - Motivation
 - New MVA optimization
 - Comparison with old one
 - MVA optimal cuts and invariant mass distributions
- Conclusions

Motivation

From Run I analysis (Phys.Rev.D90(2014), no.11 112004)

- Strong evidence for integrated CP asymmetry in the four channels
 - $A_{CP}(B^{\pm} \rightarrow K^{\pm}\pi^{+}\pi^{-}) = +0.025 \pm 0.004 \pm 0.004 \pm 0.007$, Significance: 2.8 σ
 - $A_{CP}(B^{\pm} \rightarrow K^{\pm}K^{+}K^{-}) = -0.036 \pm 0.004 \pm 0.002 \pm 0.007$, Significance: 4.3 σ

•
$$A_{CP}(B^{\pm} \rightarrow \pi^{\pm}\pi^{+}\pi^{-}) = +0.058 \pm 0.008 \pm 0.009 \pm 0.007$$
, Significance: 4.2 σ

• $A_{CP}(B^{\pm} \to \pi^{\pm} K^{+} K^{-}) = -0.123 \pm 0.017 \pm 0.012 \pm 0.007$, Significance: 5.6 σ

Large asymmetries in regions of the phase space
 Particularly large in the ππ ⇔ KK re-scattering region



Hot color: Excess of B^- Cold color: Excess of B^+

- Run I: 2011+2012
- Run II: 2015+2016 $\sim 2 \times \text{Run I}$

Motivation for MVA

- Traditionally: Cuts on each variable of a set
 - You know what you are doing
 - An event is rejected based in just one variable
 - Hard for many variables
 - Computationally expensive to automate
 - Lower perfomance
- Multivariate methods (for classification): Combine the variables into a single one
 - Many ways to combine. Find the best for your case
 - Sometimes hard to understand/follow how they are combined
 - An event "rejected" by one variable may be recovered in the combination with others
 - Smooth phase space borders
 - In general, better performance

Motivation for MVA

Some methods:

- Projective Likelihood Estimator
- Fisher discriminant
- Artificial Neural Networks
- Boosted Decision Trees BDT
- Support Vector Machine

• ...

Some packages

- TMVA: Toolkit for Multivariate Data Analysis (http://tmva.sourceforge.net)
- NeuralBayes

(https://twiki.cern.ch/twiki/bin/view/Main/NeuroBayes)

 scikit-learn: Machine Learning in Python (http://scikit-learn.org/stable)

MVA in Run I analysis(Phys.Rev.D90(2014), no.11 112004)

Note: Cuts used in Phys.Rev.Lett.112, 011801 (2014), Phys.Rev.Lett.111, 101801 (2013)

Training:

Signal: MC simulation

Background: Data in a right side invariant mass window

- Input: 23 discriminating variables
- Data pre-processing: PCA (Principal component analysis)
- MVA method: BDT
- One single optimization and cut for all four channels

Training with joint MC from all channels and background from $B
ightarrow \pi \pi \pi$



MVA in Run II analysis

- Data from 2015 and 2016. May include 2017
- Pre-selection:
 - Trigger, loose cuts (*stripping*), loose PID
- Update from Run I
 - Two more variables (25 in total)
 - Extended background region: (5400 $< B_{-}m < 6300$) MeV/c^2
 - Joint 2015+2016 data
 - Method: BDT with PCA
- New MVA optimization: Reduced set of input variables
 - Too many variables
 - Many are highly correlated
 - Some don't show good MC-data agreement
 - Set of just 10 discriminating variables
 - Method: BDT, no need for PCA

Showing some sample plots next slides

Correlations for $B \rightarrow K \pi \pi$ background 25 variables

Correlation Matrix (background)



Similar correlations hold for the signal MC and for other channels

New input variables for $B \to \pi \pi \pi$



Choice of MVA method for $B ightarrow K \pi \pi$

Background rejection versus Signal efficiency Background rejection 0.95 MVA 0.9 PChri BDT5100 0.85 BDToot BDT3100 0.8 AL P BDT BDTPCAopt 0.75 BDTPCA5100 BDTPCA3100 0.7 BDTRCA PDFRS 0.65 PDERSPCA BoostedFisher 0.6 Likelihood LikelihoodPCA 0.55 CutsGA 0.5 0.55 0.6 0.65 0.7 0.75 0.8 0.85 0.9 0.95 Signal efficiency Cut efficiencies and optimal cut value rtraining check for classifier: BDT Signal(%) for B = 10% Signal (training sample) Signal e Recknessed ettis Background () Overtraining checks train test **MLPChris** 93.2 93.5 BDT 92.7 92.8 BDTPCA 91.6 91.8 28.93 when cutting at -0.02 -0.5 -0.4 -0.3 -0.2 -0.1 0 0.1 0.2 -0.4 -0.3 0.1 0.2 0.3 0.4 0.4 -0.2

BDT response

RENAFAE Workshop July, 2018

Cut value applied on BDT output

Performances with 25 and 10 variables (loose PID requirement)



★ 3 → 3

Performances of optimizations with 25 and 10 variables

- From training/testing samples: Set with 10 variables gives performances somewhat inferior compared to the set with 25 variables
- However, when MVA cut is optimized combined, with tight PID requirements, we can easily recover the lost efficiency with a small penalty in significance
- Strategy for MVA cut optimization:
 - S is taken from MC sample: S_{MC}
 - (S + B) is the number of events in the signal region (abs(B_m - 5284) < 40) in the data spectrum

$$Significance = \frac{S_{MC}}{\sqrt{(S+B)_{Data}}}$$

 Choose MVA cut as the value close to the maximum significance that gives good signal efficiency

12/16

MVA cut optimization - $B \rightarrow K \pi \pi$ (With tight PID requirements)



Cut is a compromise between high significance and high efficiency *Common: Training with joint MC from all channels and background from pipping Same MVA cut for all channels

13/16

MVA selected events invariant mass fits





・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

14/16

MVA selected events invariant mass fits

 $B \rightarrow KK\pi$



15/16

Conclusions

Run II analysis of $B \rightarrow hhh$: MVA with 10 input variables

- Cleaner, simpler analysis
- Less correlations
- Better MC-data agreement
- MVA method: BDT (Was 25 vars with BDT+PCA)
- Final performances similar to the ones with 25 vars Both for the self or for the common optimizations
- Background levels are acceptable and well behaved