Machine learning in LHCb: triggering and flavour tagging systems

Tatiana Likhomanenko

Machine learning in LHCb: triggering and flavour tagging systems

LHC data structure



LHC event



LHC data structure



Event definition

- Proton-proton bunches collide \rangle
- Sample: one proton-proton bunches \rangle collision, called Event
- Event consists of the tracks and secondary vertices (SV), where particles are produced
- Features: a track, SV and its products \rangle physical characteristics reconstructed from the detectors (momentum, mass, angles, impact parameter)





Event processing



LHC data structure



Machine learning in LHCb: triggering and flavour tagging systems

Trigger system





What is it?

- Select events to store them for offline processing
- Should efficiently select interesting events
- Interesting event is an event that contains at least one SV where necessary particle and products are produced
- > Output rate for trigger system is limited

Trigger system



Run-II topological trigger

- HLT-1 track is looking for either one super high \rangle PT or high displacement track
- HLT-1 2-body SV classifier is looking for two \rangle tracks making a vertex
- HLT-2 improved topo classifier uses full \rangle reconstructed event to look for 2, 3, 4 and more tracks making a vertex







Machine learning problem

- Training data are set of SVs for all events \rangle
- Monte Carlo sample (used as signal-like) were simulated for various types of interesting events (different decays)
- Minimum bias data (real data for a small period of time) are used as background-like
- Output rate is fixed, thus, false positive rate (FPR) for events is fixed
- Goal is to improve efficiency for each type of signal events.



If at least one SVR in event passed all stages, whole event passes trigger

Event representation



How to measure quality?



Trigger system



ROC curve, computed for events

- Optimize true positive rate (TPR) for fixed FPR for events
- > Weight signal events in such way that channels have the same amount of events.
- Optimize ROC curve in a small FPR region



ROC curve interpretation

11

Machine learning in LHCb: triggering and flavour tagging systems

Trigger system: random forest trick

Simulated signal event contains at least one interesting SV, but not each SV should be interesting

Random forest for SVRs selection

- Train random forest (RF) on SVRs \rangle
- > Select top-1, top-2 SVs by RF predictions for each signal event
- > Train classifier on selected SVs

Trigger system: random forest trick





Machine learning in LHCb: triggering and flavour tagging systems

Trigger system: real-time



Online processing

There are two possibilities to speed up prediction operation:

- Bonsai boosted decision tree format (BBDT) \rangle
- Post-pruning \rangle



What is MatrixNet

- > Yandex machine leaning algorithm
- Gradient Boosting over oblivious Decision Trees \rangle
- > Feature binarization (like feature hashing)
- Classification, Regression, Ranking \rangle

Real-time trigger system

17

BBDT

- Features hashing using bins before training
- Converting decision trees to \rangle n-dimensional table (lookup table)
- Table size is limited in RAM (1Gb), thus count of bins for \rangle each features should be small (5 bins for each of 12 features)
- Discretization reduces the quality
- Prediction operation takes one reading from the table







BBDT, results



T MN special	 rate: 2.5 kHz
e MN	 rate: 4. kHz



Post-pruning

- Train MatrixNet (MN) with several thousands trees \rangle
- Reduce this amount of trees to a hundred
- Greedily choose trees in a sequence from the initial ensemble to \rangle minimize a modified loss function:

 $\sum_{sianal} \log\left(1 + e^{-F(x)}\right) + \sum_{backaround} e^{F(x)}$ background signal

At the same time change values in leaves (tree structure is preserved) \rangle



Post-pruning, results





Trigger results



https://github.com/yandexdataschool/LHCb-topo-trigger



Machine learning in LHCb: triggering and flavour tagging systems

sPlot technique





Solution for what?

- Monte Carlo is not-well simulated
- Need to work with real unlabeled data
- \rangle background data
- \rangle physics) the mass pdfs for signal and background.
- How to restore signal/bck pdfs for other features? >

Need somehow to label real data: want to restore for features their distributions for the signal and

Our main knowledge is the mass distribution for real data from which we can extract (using some



Feature initial distributions



sPlot technique



Two mass bins

Proportion of events inside bins



 $Bin1: w_{b_1}f_b + w_{s_1}f_s$ $Bin2: w_{b_2}f_b + w_{s_2}f_s$

sPlot technique



 $*w_{b_2} + (-w_{b_1})$

will obtain initial signal distribution



Reconstruction



sPlot technique





More bins: sWeight

- > Equivalent to some optimization problem
- > Have explicit solution
- > Produce weights (sWeight) for each event
- > Feature pdf with sWeight will be signal pdf

sPlot technique

Proportion of events inside bins



How to reweight?



Machine learning in LHCb: triggering and flavour tagging systems

Tagging system (not official)



What is it?

- Event has a signal decay part \rangle
- The signal decay part can be produced from b quark or anti-b quark
- > The system should effectively predict the source of the signal decay (b quark or anti-b quark)
- An intermediate B-meson in the signal decay part can oscillate \rangle
- The tagging system prediction P(anti-b quark) will allow to >measure the oscillation effects



Tagging particles (goal)





Tagging particles (training)





Confidence interval to asymmetry

- Construct probabilities for $b \rightarrow$ final state and anti-b \rightarrow final state using \rangle P(anti-b quark)
- Take relation of probabilities, called asymmetry \rangle
- Confidence interval for parameter of interest $(A_m measured)$: \rangle $\frac{\sqrt{1-A_m^2}}{\sqrt{\varepsilon_{tag}N}(1-2\omega)}$

$$\sigma_A = \frac{\sigma_{A_m}}{1 - 2\omega} = \frac{1}{\sqrt{2\omega}}$$

Tagging system should maximize effective efficiency (ω - mistag probability):

$$\varepsilon_{eff} = \varepsilon_{tag} D^2 = 0$$

 $\varepsilon_{tag}(1-2\omega)^2$



Current tagging system

- Take $B^{\pm} \rightarrow J/\psi K^{\pm}$
- Apply sPlot technique to real data to extract signal-like data (sWeights)
- Choose one tagging track (physical selections, like PID and max PT track) for each event >
- Train exclusive taggers for SS (same side) tagging particles and OS (opposite side) \rangle
- Target for classifier is 'right tagged' label >
- Combine all taggers to one (probabilistic model) to obtain P(anti-b quark)



Inclusive tagging system (new ideas)

- Don't apply physical selections, except loose cuts on PID and \rangle track ghost
- > For each event use full topological information: all tracks and possible vertices
- Use probabilistic model for tracks and vertices to obtain \rangle P(anti-b quark)
- Maximize ROC AUC score



Restrictions

- > Maximize effective efficiency
- > Necessary to calibrate output to probability
- > Distribution should be symmetry for b-quark (B⁻) and anti-b quark (B⁺)
- Flatness for B-mass, lifetime, momentum (to simplify further analysis)



Track-based inclusive tagging









Track-based inclusive tagging

- Target is BSign*trackSign > 0 (avoid to define tagging particle)
- Classifier will return P(track has same sign as B | B sign) \rangle
- Suppose:

P(track has same sign as B | B sign) =

P(B has same sign as track | track sign)

- AUC 0.5134
- Calibrate output to probability (isotonic calibration)





Vertex-based inclusive tagging

- Target is BSign*svSign > 0 \rangle
- Classifier will return P(vertex has same sign as B | B sign)
- Suppose:

P(vertex has same sign as B | B sign) =

P(B has same sign as vertex | vertex sign)

- AUC 0.5544
- Calibrate output to probability (logistic calibration)





Calibration

- Platt's calibration: logistic regression over the classifier output
- > Bin method:
 - #(same sign in bin) / #(in bin)
 - fit by linear function
- > Isotonic regression: monotonic function (extend the bin method)





Isotonic calibration







Probabilistic model for events



where



Tagging system

$$P(\text{track/vertex}|B^+) = \alpha$$

$$P(\text{track/vertex}|B^-) = (1)$$

$$= \frac{\alpha}{1 + \alpha}, \qquad [1]$$

 $p_{mistag} = min(p(B^+), p(B^-))$



anti-b quark probability

- > Isotonic calibration
- > Produced tied probabilities
- > Add small normal noise

0.001 * normal(0, 1)







Check B symmetry (before calibration)

- \rangle



Check B symmetry (after calibration)

- KS: 0.01683



ROC for events (not official)

- AUC score 0.64 (for current tagging system 0.566) \rangle
- What about calibration?







Check calibration (with B-symmetry)







Thanks for attention

Contacts

Likhomanenko Tatiana researcher-developer



<u>antares@yandex-team.ru</u>

References

- <u>record/1384380/files/LHCb-PUB-2011-016.pdf</u>, <u>http://arxiv.org/abs/1510.00572</u>
- <u>1381330/files/CERN-THESIS-2015-040.pdf</u>
- http://arxiv.org/pdf/1211.0025.pdf, http://arxiv.org/pdf/1511.00213
- sPlot technique: http://arxiv.org/pdf/physics/0402083

Tagging system

Trigger system: <u>https://github.com/yandexdataschool/LHCb-topo-trigger</u>, <u>https://cdsweb.cern.ch/</u>

Tagging system: <u>https://github.com/tata-antares/tagging_LHCb</u>, <u>https://inspirehep.net/record/</u>

Calibration: <u>http://fastml.com/classifier-calibration-with-platts-scaling-and-isotonic-regression/</u>,



- After the mass maximum likelihood fitting we know for each event y $p_s(y), p_b(y)$
- which are probabilities of an event to be signal and background.
- Will reconstruct number of *signal* events in a *particular* histogram bin for the reconstructed feature. Introduce unknown probability that the signal/bck event will be in a particular bin:
- The total amount of signal/bck events obtained from the fit:
- p_s, p_b
 - N_s, N_b

The random variable, number of signal events in the bin:

where $w_s(y)$ are sPlot weights and are a subject to find.

Property: an estimation should be unbiased:

Corollary:

 $p_s N_s = \mathbb{E}X =$

$$=\sum_{y} w_s(y)($$

- $X = \sum w_s(y) \mathbb{I}_{y \in bin}$

 - $\mathbb{E}X = p_s N_s$

$$\sum_{y} w_s(y) \mathbb{E}\mathbb{I}_{y \in bin} =$$

 $(p_s p_s(y) + p_b p_b(y))$

Since the previous equation should hold for all possible p_s, p_b , we get two equalities:



Then we can guarantee that mean input of background are 0 (the expectation is zero, but observed number will not be zero due to the deviation)

$$\sum w_s(y) p_s p_s(y)$$

y

 $0 = \sum w_s(y) p_b p_b(y)$

$$\sum_{y} w_s(y) p_s(y)$$

$$\bigg] w_s(y) p_b(y)$$

Assumption of the linearity:

Then:

$$\begin{pmatrix} V_{bb} & V_{bs} \\ V_{sb} & V_{ss} \end{pmatrix} * \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} 0 \\ N_s \end{pmatrix}$$
$$j = \sum p_i(y) * p_j(y)$$

where
$$V_{ij} = \sum_{y} p_i(y) * p_j(y)$$

The assumption of the linearity is correct because apart from having correct mean, we should also minimize the variation of the reconstructed variable.

$$\mathbb{V}X = \sum_{y} w_s^2(y) \mathbb{V}\mathbb{I}_{y \in bin} \le \sum_{y} w_s^2(y)$$

 $w_s(y) = a_1 p_b(y) + a_2 p_s(y)$

Minimization problem



Lagrangian:

$$\mathcal{L} = \sum_{y} w_s^2(y) + \lambda_1 \left(\sum_{y} w_s(y) p_b(y) \right) + \lambda_2 \left(\sum_{y} w_s(y) p_s(y) - N_s \right)$$
$$0 = \frac{\partial \mathcal{L}}{\partial w_s(y)} = 2w_s(y) + \lambda_1 p_b(y) + \lambda_2 p_s(y)$$

It holds for each event, thus we are getting needed the linear dependence

 $\sum w_s^2(y) \to min$

 $\sum w_s(y)p_b(y) = 0$

 $\sum w_s(y)p_s(y) = N_s$

Main assumption

The mass must be uncorrela $p_s(mass, feature) = p_s(y)$ Then it holds for a particular $n_s(mass, bin) = n$

- The mass must be uncorrelated with the reconstructed feature:
- $p_s(mass, feature) = p_s(y, feature) = p_s(y)p_s(feature)$
 - Then it holds for a particular bin for the reconstructed feature:
 - $p_s(mass, bin) = p_s(y, feature) = p_s(y)p_s$