

# Ensemble Learning in HEP

Kateřina Hladká

CTU in Prague

June 23, 2019

Supervisor: Ing. V. Kůs, Ph.D.

# Outline

- 1  $D^0$  Decay
- 2 Data sets
- 3 Features
- 4 Data Pre-processing
- 5 Ensemble methods
  - Random Forests
  - AdaBoost
- 6 Binary Classification Metrics and Criteria of Optimization
- 7 Results for Random Forests
- 8 Results for AdaBoost
- 9 Significance Estimations

# $D^0$ Decay

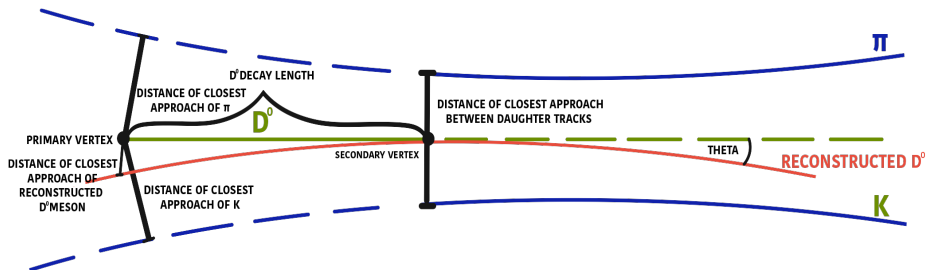


Figure:  $D^0$  Decay on STAR Experiment.

# Data sets

- Monte Carlo signal simulation data set
- Real data
  - ▶ Unlike-sign pairs
  - ▶ Like-sign pairs

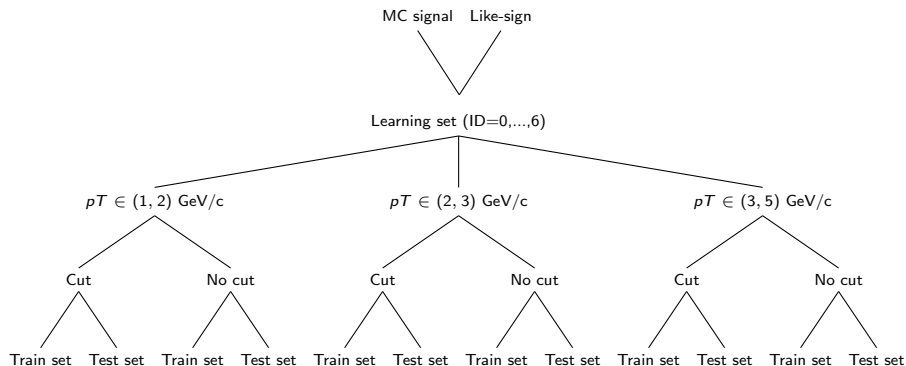
**Goal:** Determine, which unlike-sign pairs were created during  $D^0$  decay process (signal) and which were not (background).

# Features

ID	Feature Name	Cut	Unit
0	pi1_dca	$0.002 < \text{pi1\_dca} < 0.2$	$[\mu\text{m}]$
1	k_dca	$0.002 < \text{k\_dca} < 0.2$	$[\mu\text{m}]$
2	D_decayL	$0.0005 < \text{D\_decayL} < 0.2$	$[\mu\text{m}]$
3	dcaDaughters	$\text{dcaDaughters} < 0.02$	$[\mu\text{m}]$
4	cosTheta	$\text{cosTheta} > 0.7$	
5	dcaD0ToPv	$\text{dcaD0ToPv} < 0.05$	$[\mu\text{m}]$
6	D_cosThetaStar	$-1 \leq \text{D\_cosThetaStar} \leq 1$	
7	D_mass	not part of learning set	$[\text{GeV}/c^2]$
8	pT	not part of learning set	$[\text{GeV}/c]$

Table: Features for  $D^0$  classification task.

# Data Pre-processing



# Ensemble methods

- Random Forests
- AdaBoost

# Decision Trees and Random Forests

- Goal:  $\Omega = \hat{\Omega}_0 \cup \hat{\Omega}_1 = \hat{\Omega}_{\text{node } 1} \cup \dots \cup \hat{\Omega}_{\text{node } J}$ , where  $\{\text{node } 1, \dots, \text{node } J\}$  is set of all nodes.
- Strategy:
  - ▶ Find the best split of the parent node  $u$  to obtain left and right child nodes with purer subsets
  - ▶ Split nodes until stopping criterion is met

## Definition

*Gini impurity measure of node  $u$  is*

$$i_G(u) = \sum_{k=0}^1 p_u^k (1 - p_u^k), \quad (1)$$

*where  $p_u^k$  is ratio of samples of class  $k$  in the node  $u$  and all samples in the node  $u$ ,  $k = 0, 1$ .*



- Ensemble of  $M$  trees
- $M$  new training sets created using Bootstrap Aggregating

## Definition

Let  $\alpha = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} = (\mathbf{X}, \mathbf{y})$  be a training set,  $M \in \mathbb{N}$ . Then bootstrap aggregating means generating  $M$  new training sets  $\alpha_1, \dots, \alpha_M$ , where for  $m \in \{1, \dots, M\}$   $\alpha_m = \{(\mathbf{x}_{j_1}, y_{j_1}), \dots, (\mathbf{x}_{j_n}, y_{j_n})\}$  is created by sampling from  $\alpha$  uniformly and with replacement.

- In each node  $K$  features are randomly selected
- The best split is found among those  $K$  features, not among all of them
- Masking reduction

## Definition

Let  $\varphi_1, \dots, \varphi_M$  be a set of trees of a random forest. Let  $(X_1, \dots, X_r)$  be a vector of features. Then importance of feature  $X_j$ ,  $j \in \{1, \dots, r\}$ , is defined as

$$\text{Imp}(X_j) = \frac{1}{M} \sum_{m=1}^M \sum_{u \in \varphi_m} 1_{j_u=j} \left[ \frac{N_u}{N} \Delta i(d_u^*, u) \right], \quad (2)$$

where  $N_u$  denotes number of samples in node  $u$ ,  $N$  denotes number of all samples in data set,  $j_u$  denotes the identifier of the feature used to split node  $u$  and  $d_u^*$  is optimal binary split of the node  $u$ .

# AdaBoost

- Ensemble of weak learners  $\varphi_1, \dots, \varphi_A$
- Let  $\omega \in \mathbb{R}^q$ , where  $q = |\alpha_{train}|$  and  $[\omega] = (\frac{1}{q}, \dots, \frac{1}{q})$ .

$$Er = \sum_{(\mathbf{x}_i, y_i) \in \alpha_{train}} [\omega]_i \mathbf{1}_{\varphi(\mathbf{x}_i=y_i)} \quad (3)$$

$$V = \frac{1}{2} \ln \left( \frac{1 - Er}{Er} \right) \quad (4)$$

- Building ensemble of  $A$  estimators, where

$$[\omega_a]_i = \frac{[\omega_{a-1}]_i}{c_a} \times \begin{cases} e^{-V_a}, & \text{if } \varphi_{a-1}(\mathbf{x}_i) = y_i, \\ e^{V_a}, & \text{if } \varphi_{a-1}(\mathbf{x}_i) \neq y_i, \end{cases} \quad (5)$$

- Final decision

$$\varphi_{AdaBoost}(\mathbf{x}) = \begin{cases} 1, & \text{if } \sum_{a=1}^A V_a \cdot \mathbf{1}_{\varphi_a(\mathbf{x})=1} \geq \sum_{a=1}^A V_a \cdot \mathbf{1}_{\varphi_a(\mathbf{x})=0}, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

# Binary Classification Metrics and Criteria of Optimization

## Definition

$$\text{Precision} = \frac{TP}{TP + FP} \quad (7)$$

$$\text{Recall} = \text{TPR} = \frac{TP}{TP + FN} \quad (8)$$

$$\text{FPR} = \frac{FP}{FP + TN} \quad (9)$$

$$\text{ACC} = \frac{TP + TN}{TP + TN + FN + FP} \quad (10)$$

$$\text{Signifikance} = \frac{TP}{\sqrt{TP + FP}} \quad (11)$$

$$\text{F1 score} = 2 \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12)$$

# Results for Random Forests

	$pT \in (1, 2)$	$pT \in (2, 3)$	$pT \in (3, 5)$
maximum depth	8/8	7/7	7/7
ACC (train set)	0,9214/ <b>0,9524</b>	0,8752/ <b>0,9028</b>	0,9318/ <b>0,9133</b>
ACC (test set)	0,9211/ <b>0,9524</b>	0,8728/ <b>0,9020</b>	0,9311/ <b>0,9135</b>
Precision (train set)	0,7912/ <b>0,9580</b>	0,8796/ <b>0,9466</b>	0,9501/ <b>0,9322</b>
Precision (test set)	0,7904/ <b>0,9564</b>	0,8766/ <b>0,9466</b>	0,9504/ <b>0,9326</b>
Recall (train set)	0,6165/ <b>0,5579</b>	0,8150/ <b>0,7385</b>	0,9757/ <b>0,9684</b>
Recall (test set)	0,6165/ <b>0,5578</b>	0,8126/ <b>0,7360</b>	0,9747/ <b>0,9683</b>
F1 score (train set)	0,6930/ <b>0,7052</b>	0,8460/ <b>0,8297</b>	0,9627/ <b>0,9499</b>
F1 score (test set)	0,6927/ <b>0,7046</b>	0,8434/ <b>0,8281</b>	0,9624/ <b>0,9501</b>
# signal (train set)	214370/ <b>368801</b>	263017/ <b>447439</b>	615599/ <b>1055522</b>
# background (train set)	1274323/ <b>13247861</b>	362241/ <b>947728</b>	66182/ <b>186852</b>
# signal (test set)	92023/ <b>157794</b>	112959/ <b>191877</b>	264274/ <b>452674</b>
# background (test set)	545989/ <b>1392205</b>	155010/ <b>406052</b>	27918/ <b>79773</b>

**Table:** Random forests with cuts/**without cuts** and optimal threshold  $\delta^*$  with respect to maximal F1 score.

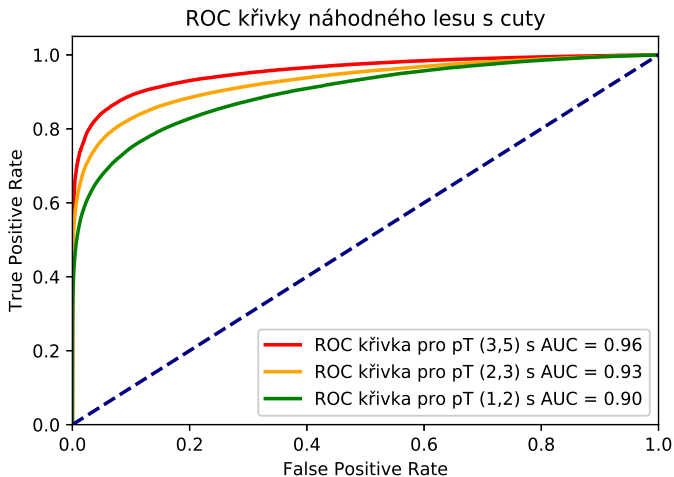


Figure: ROC curves of random forests with cuts.

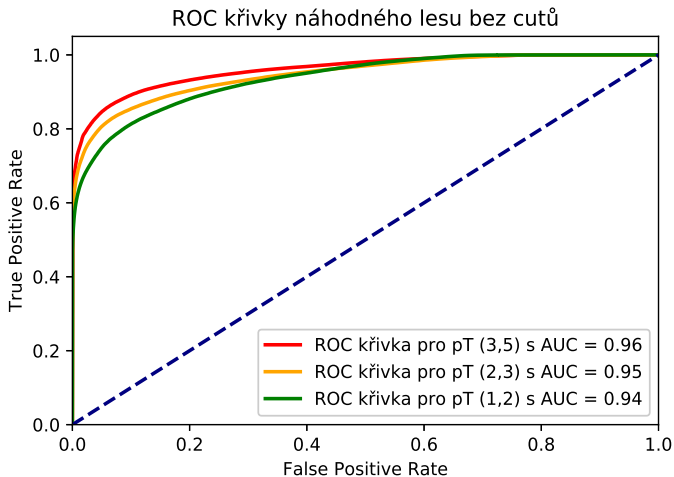
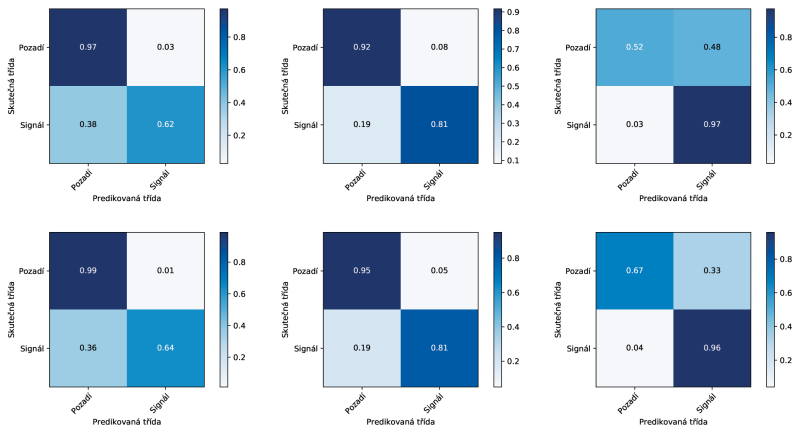


Figure: ROC curves of random forests without cuts.



**Figure:** First line: confusion matrices for  $pT \in (1, 2)$ ,  $pT \in (2, 3)$ ,  $pT \in (3, 5)$ .  
 Second line: confusion matrices for  $pT \in (1, 2)$ ,  $pT \in (2, 3)$ ,  $pT \in (3, 5)$  **no cuts**.

	$pT \in (1, 2)$	$pT \in (2, 3)$	$pT \in (3, 5)$
# unlike-sign pairs	1846337/ <b>4713795</b>	530396/ <b>1388118</b>	97953/ <b>278494</b>
signal/background ratio	0,014/ <b>0,006</b>	0,044/ <b>0,018</b>	0,329/ <b>0,261</b>

**Table:** Results for random forest applied to real collision, trained on cutted/**not cutted** data.



# Results for AdaBoost

	Bin 1	Bin 2	Bin 3
ACC (train set)	0,9184/ <b>0,9512</b>	0,8678/ <b>0,9038</b>	0,9268/ <b>0,9143</b>
ACC (test set)	0,9185/ <b>0,9513</b>	0,8663/ <b>0,9027</b>	0,9266/ <b>0,9149</b>
Precision (train set)	0,8541/ <b>0,9281</b>	0,8813/ <b>0,9061</b>	0,9481/ <b>0,9382</b>
Precision (test set) $y$	0,8547/ <b>0,9273</b>	0,8798/ <b>0,9054</b>	0,9486/ <b>0,9385</b>
Recall (train set)	0,5223/ <b>0,5651</b>	0,7926/ <b>0,7809</b>	0,9722/ <b>0,9625</b>
Recall (test set)	0,5243/ <b>0,5659</b>	0,7909/ <b>0,7782</b>	0,9715/ <b>0,9630</b>
F1 score (train set)	0,6482/ <b>0,7025</b>	0,8346/ <b>0,8388</b>	0,9600/ <b>0,9502</b>
F1 score (test set)	0,6499/ <b>0,7028</b>	0,8330/ <b>0,8370</b>	0,9599/ <b>0,9506</b>
# signal (train set)	214370/ <b>368801</b>	263017/ <b>447439</b>	615599/ <b>1055522</b>
# background (train set)	1274323/ <b>3247861</b>	362241/ <b>947728</b>	66182/ <b>186852</b>
# signal (test set)	92023/ <b>157794</b>	112959/ <b>191877</b>	264274/ <b>452674</b>
# background (test set)	545989/ <b>1392205</b>	155010/ <b>406052</b>	27918/ <b>79773</b>

Table: AdaBoost with/**without** cuts and optimal threshold  $\delta^*$  with respect to maximal F1 score.

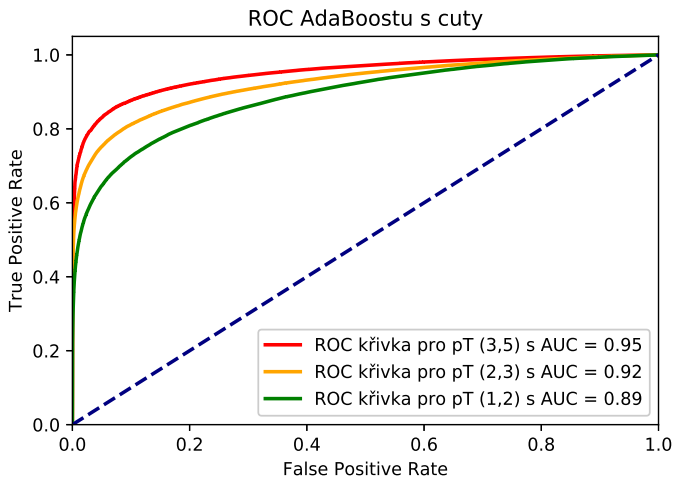


Figure: ROC curves of AdaBoost with cuts.

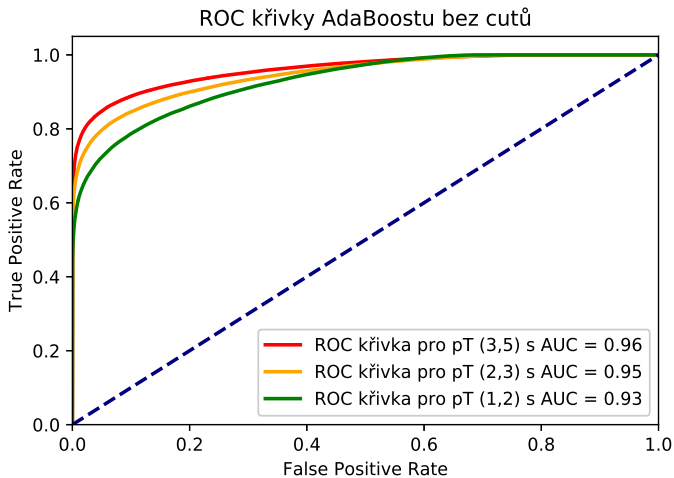


Figure: ROC curves of AdaBoost without cuts.

	$pT \in (1, 2)$	$pT \in (2, 3)$	$pT \in (3, 5)$
# unlike-sign pairs	1846337/ <b>4713795</b>	530396/ <b>1388118</b>	97953/ <b>278494</b>
signal/background ratio	0,008/ <b>0,003</b>	0,042/ <b>0,020</b>	0,345/ <b>0,227</b>

**Table:** Results for AdaBoost applied to real collision, trained on cutted/**not cutted** data.

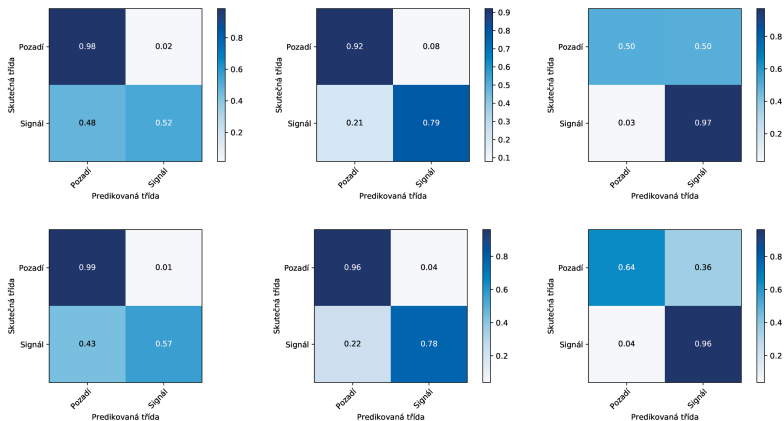


Figure: First line: confusion matrices for  $pT \in (1, 2)$ ,  $pT \in (2, 3)$ ,  $pT \in (3, 5)$ .  
 Second line: confusion matrices for  $pT \in (1, 2)$ ,  $pT \in (2, 3)$ ,  $pT \in (3, 5)$  **no cuts**.

# Significance Estimations (Sneak Peek)

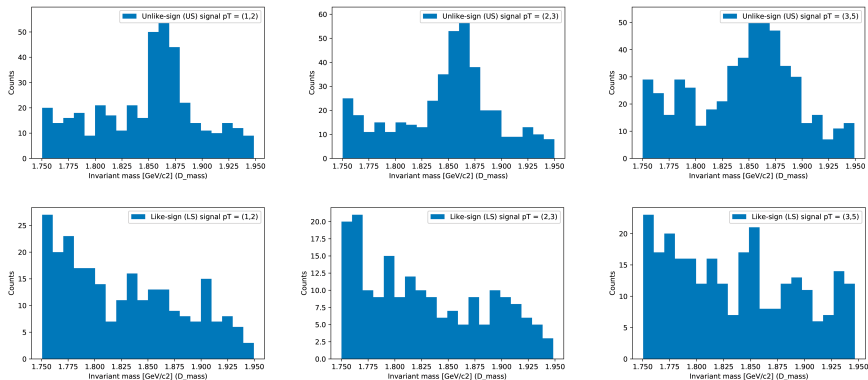


Figure: Unlike-sign and like-sign histograms.

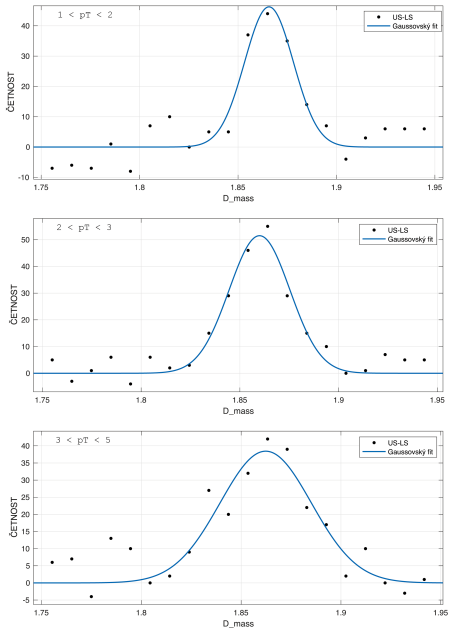







Figure: Future goal.

# References

-  WOLFRAM FISCHER *Run overview of the Relativistic Heavy Ion Collider*. [online, cit. 2019-03-30], <http://www.rhichome.bnl.gov/RHIC/Runs/>
-  D. TLUSTÝ *A Study of Open Charm Production in  $p+p$  Collisions at STAR*. Dizertační práce, ČVUT, FJFI
-  GILLES LOUPPE *Understanding random forests from theory to practise*. Dizertační práce, University of Liège, Faculty of Applied Sciences, Department of Electrical Engineering and Computer Science, 2014
-  L. BRIEMAN *Random Forests*, Statistics Department, University of California, 2001
-  Y. FREUND, R. E. SCHAPIRE *Boosting, Foundations and Algorithms*. The MIT Press, Cambridge, Massachusetts, London, England, 2012.



Thank You for Your attention.