

“Flavours of Physics” challenge: 2nd place solution

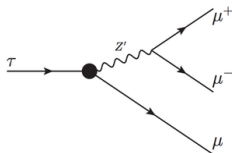
Alexander V. Gramolin
(Budker Institute of Nuclear Physics, Novosibirsk, Russia)



ALEPH Workshop @ NIPS 2015
Montreal, Canada
December 11, 2015



“Flavours of Physics” machine learning challenge



- Three types of events:

Signal
(simulated):

$$\tau^- \rightarrow \mu^- \mu^+ \mu^-$$
$$m_\tau = 1776.82 \text{ MeV}$$

Control
(simulated + real):

$$D_s^- \rightarrow \phi(\mu^+ \mu^-)\pi^-$$
$$m_{D_s} = 1968.30 \text{ MeV}$$

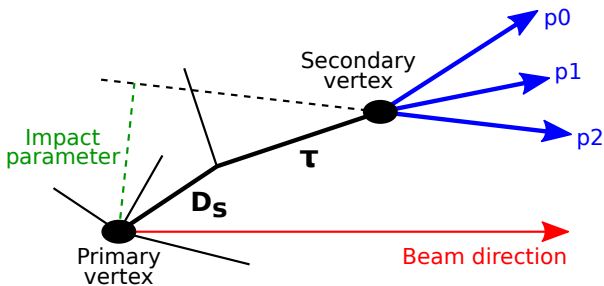
Background
(real):

$$D_s^- \rightarrow \eta(\mu^+ \mu^- \gamma)\mu^- \bar{\nu}_\mu$$
$$D_{(s)}^- \rightarrow K^+ \pi^- \pi^-$$
$$D_{(s)}^- \rightarrow \pi^+ \pi^- \pi^-$$

- Four datasets were provided:
 - **training.csv** (67 553 labeled events, 41 674 signal + 25 879 background)
 - **test.csv** (855 819 unlabeled events, signal + control + background)
 - **check_agreement.csv** (331 147 control events)
 - to check that the classifier is not exploiting the imperfections of simulation
 - **check_correlation.csv** (5 514 background events)
 - to check that the prediction is not too correlated with the τ mass

Original features

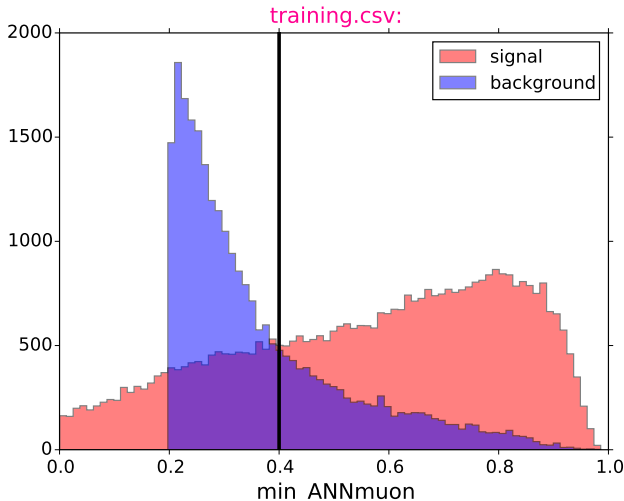
- 46 primary features:
 - 10 kinematic features (momenta, transverse momenta, and pseudorapidities of p_0 , p_1 , and p_2 ; transverse momentum of the mother particle)
 - 35 "geometric" features (impact parameters, track isolation variables, flight distance and lifetime of the mother particle, etc.)
 - SPDhits (number of hits in the Scintillating Pad Detector, SPD)



- 3 auxiliary features (absent in the test dataset):
 - production (τ production mechanism: $D_s^- \rightarrow \tau$, $D^- \rightarrow \tau$, or $X_b \rightarrow \tau$)
 - min_ANNmuon (muon identification variable)
 - mass (mass of the mother particle)

The “min_ANNmuon” feature

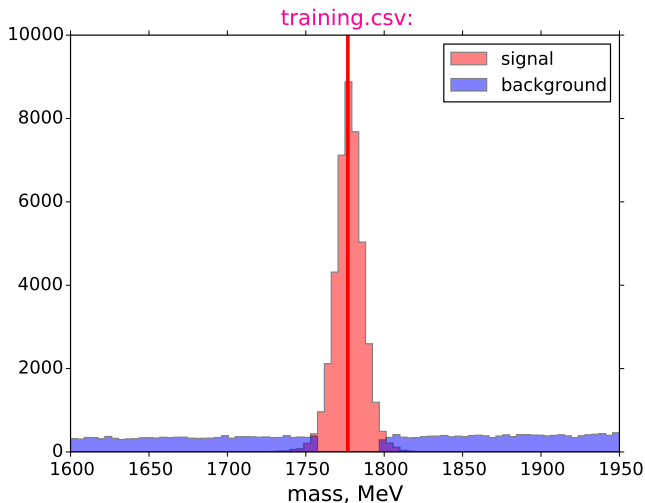
Muon identification variable obtained using an artificial neural network (ANN)



Only events with $\text{min_ANNmuon} > 0.4$ are used to score a model
⇒ use only these events for training (30 213 signal + 7 799 background events)

The “mass” feature

The real events collected by LHCb were used as background. To make sure that there are no $\tau \rightarrow 3\mu$ decays among them, events were removed from the τ mass window.



⇒ mass is a “golden feature” that can separate signal and background almost perfectly!

How to reconstruct the mass: the first way

1. Calculate the longitudinal momenta of p_0 , p_1 , and p_2 :

$$\mathbf{p}_{0,z} = p_{0,pt} \cdot \sinh(p_{0,\eta}), \quad \mathbf{p}_{1,z} = p_{1,pt} \cdot \sinh(p_{1,\eta}), \quad \mathbf{p}_{2,z} = p_{2,pt} \cdot \sinh(p_{2,\eta}).$$

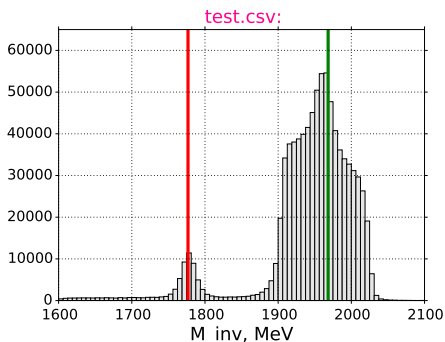
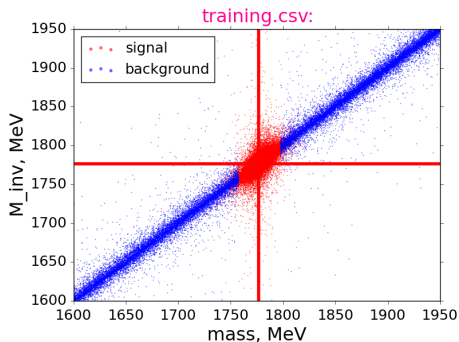
2. Assume that p_0 , p_1 , and p_2 are muons and calculate their full energies:

$$\mathbf{E}_0 = \sqrt{m_\mu^2 + p_{0,p}^2}, \quad \mathbf{E}_1 = \sqrt{m_\mu^2 + p_{1,p}^2}, \quad \mathbf{E}_2 = \sqrt{m_\mu^2 + p_{2,p}^2}, \quad m_\mu = 105.66 \text{ MeV}.$$

3. Calculate the kinematic parameters of the mother particle:

$$\mathbf{E} = \mathbf{E}_0 + \mathbf{E}_1 + \mathbf{E}_2, \quad \mathbf{p}_z = p_{0,z} + p_{1,z} + p_{2,z},$$

$$\mathbf{M}_{inv} = \sqrt{\mathbf{E}^2 - (\mathbf{p}_t^2 + \mathbf{p}_z^2)}.$$



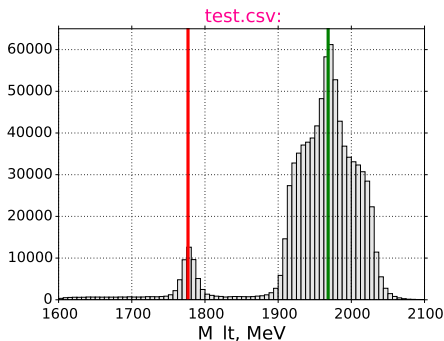
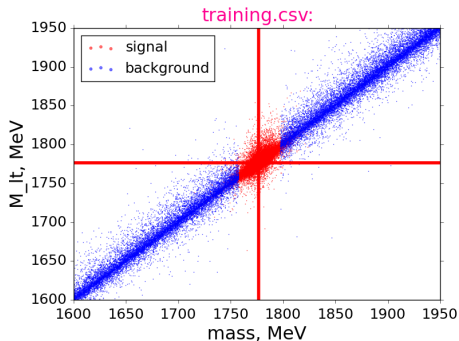
How to reconstruct the mass: the second way

Alternatively, the mass of the mother particle can be calculated from **LifeTime** and **FlightDistance** using special relativity:

$$\beta\gamma \approx \frac{\text{FlightDistance}}{\text{LifeTime} \cdot c}, \quad \text{where } \beta = \frac{v}{c}, \quad \gamma = \frac{1}{\sqrt{1 - \beta^2}},$$

and $c = 299.792458 \frac{\text{m}}{\mu\text{s}}$ is the speed of light. Then

$$\mathbf{M_It} = \frac{\sqrt{pt^2 + pz^2}}{\beta\gamma}.$$



⇒ **M_It** is a more accurate estimate of the mass than **M_inv**!

Additional kinematic features

- Features related to the mother particle:

$$\text{flag_M} = \begin{cases} 1 & \text{if } |M_{\text{lt}} - m_{\tau} - 1.44| < 17, \\ 0 & \text{otherwise,} \end{cases}$$

$$\text{Delta_M} = M_{\text{lt}} - M_{\text{inv}},$$

$$\text{Delta_E} = \sqrt{M_{\text{lt}}^2 + (p_t^2 + p_z^2)} - E,$$

$$\text{gamma} = \frac{E}{M_{\text{inv}}}, \quad \text{beta} = \frac{\sqrt{\text{gamma}^2 - 1}}{\text{gamma}}.$$

- Features related to the final-state particles p0, p1, and p2:

$$\text{E0_ratio} = \frac{E_0}{E}, \quad \text{E1_ratio} = \frac{E_1}{E}, \quad \text{E2_ratio} = \frac{E_2}{E},$$

$$\text{p0_pt_ratio} = \frac{p_0\text{-pt}}{p_t}, \quad \text{p1_pt_ratio} = \frac{p_1\text{-pt}}{p_t}, \quad \text{p2_pt_ratio} = \frac{p_2\text{-pt}}{p_t},$$

$$\text{t_coll} = \frac{p_0\text{-pt} + p_1\text{-pt} + p_2\text{-pt}}{p_t},$$

$$\text{eta_01} = p_0\text{-eta} - p_1\text{-eta}, \quad \text{eta_02} = p_0\text{-eta} - p_2\text{-eta},$$

$$\text{eta_12} = p_1\text{-eta} - p_2\text{-eta}.$$

Additional “geometric” features

- Significance of the flight distance:

$$\mathbf{FlightDistanceSig} = \frac{\mathbf{FlightDistance}}{\mathbf{FlightDistanceError}}.$$

- Sums of the original features:

$$\mathbf{DOCA_sum} = \mathbf{DOCAone} + \mathbf{DOCAtwo} + \mathbf{DOCAthree},$$

$$\mathbf{isolation_sum} = \mathbf{isolationa} + \mathbf{isolationb} + \mathbf{isolationc} \\ + \mathbf{isolationd} + \mathbf{isolatione} + \mathbf{isolationf},$$

$$\mathbf{IsoBDT_sum} = \mathbf{p0_IsoBDT} + \mathbf{p1_IsoBDT} + \mathbf{p2_IsoBDT},$$

$$\mathbf{IP_sum} = \mathbf{p0_IP} + \mathbf{p1_IP} + \mathbf{p2_IP},$$

$$\mathbf{IPSig_sum} = \mathbf{p0_IPSig} + \mathbf{p1_IPSig} + \mathbf{p2_IPSig},$$

$$\mathbf{CDF_sum} = \mathbf{CDF1} + \mathbf{CDF2} + \mathbf{CDF3}.$$

- Quality of the p0, p1, and p2 tracks:

$$\mathbf{track_Chi2Dof} = \left[(\mathbf{p0_track_Chi2Dof} - 1)^2 + (\mathbf{p1_track_Chi2Dof} - 1)^2 \right. \\ \left. + (\mathbf{p2_track_Chi2Dof} - 1)^2 \right]^{1/2}.$$

Description of the solution

- Boosted decision trees (BDTs) is a trusted and widely used tool in HEP.
- *dmlc* **XGBoost** is a popular BDT library successfully used by many Kagglers. It has recently won the “HEP meets ML” award from CERN.
- The model is based on two BDT classifiers using different sets of features and trained independently. The first classifier is **geometric** and the second one is **kinematic**. Final prediction is

$$\text{prediction} = \frac{1}{2} [w_1 \cdot \text{prediction}_1 + (1 - w_1) \cdot \text{prediction}_2], \quad w_1 = 0.78.$$

- The geometric classifier is weaker, but helps to pass the **correlation test**. To pass the **agreement test**, it was enough to not use the **SPDhits** feature.
- Source code in Python: <https://github.com/gramolin/flavours-of-physics>. It contains the following files:
 - **features.py** (implementation of the additional features)
 - **parameters.py** (tunable parameters of the model)
 - **train.py** (training the classifiers)
 - **predict.py** (making a prediction)

```
# Random seed:
```

```
random_state = 1
```

```
# Weight for the first classifier:
```

```
w1 = 0.78
```

```
# Numbers of trees:
```

```
num_trees1 = 200 # Classifier 1
```

```
num_trees2 = 100 # Classifier 2
```

```
# Parameters of the classifiers:
```

```
params = {'objective': 'binary:logistic',  
          'eta': 0.05,  
          'max_depth': 4,  
          'scale_pos_weight': 5.,  
          'silent': 1,  
          'seed': random_state}
```

```
import pandas, xgboost, features, parameters

# Read the training dataset:
train = pandas.read_csv('data/training.csv', index_col='id')
train = train[train['min_ANNmuon'] > 0.4]

# Add extra features:
train = features.add_features(train)

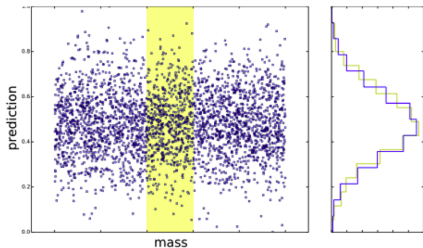
# Train the first (geometric) XGBoost classifier:
bst1 = xgboost.train(parameters.params,
                    xgboost.DMatrix(train[features.list1],
                                     train['signal']), parameters.num_trees1)
bst1.save_model('bst1.model')

# Train the second (kinematic) XGBoost classifier:
bst2 = xgboost.train(parameters.params,
                    xgboost.DMatrix(train[features.list2],
                                     train['signal']), parameters.num_trees2)
bst2.save_model('bst2.model')
```

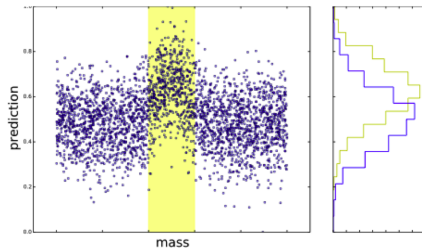
Correlation test

A Cramer–von Mises test is performed to compare the predicted values for the entire dataset with the predicted values within a rolling window.

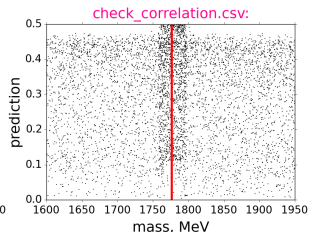
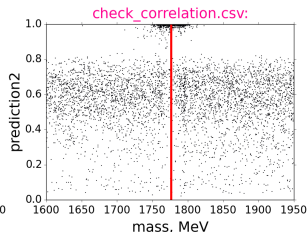
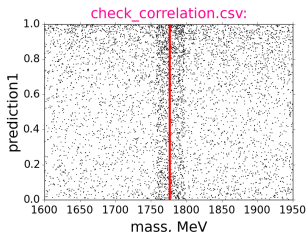
No correlation with mass:



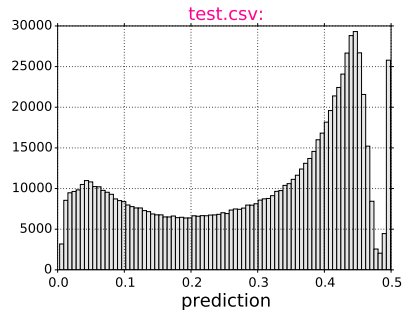
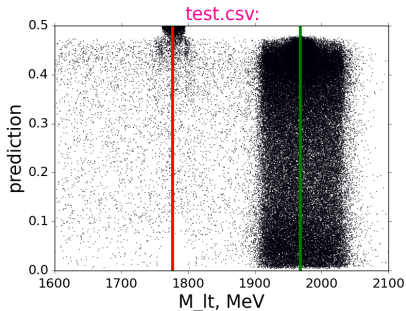
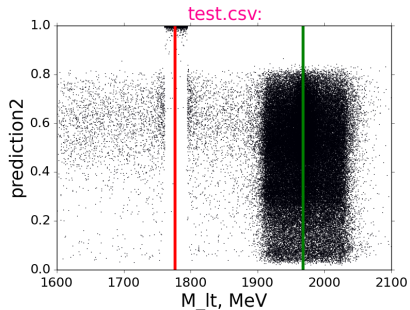
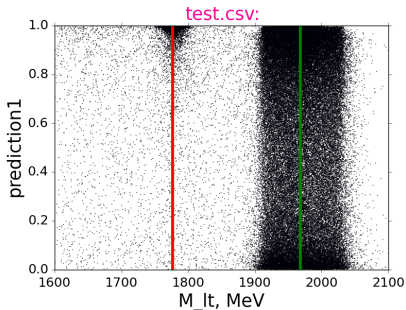
Strong correlation:



Predictions for the correlation dataset:



Predictions for the test dataset



Conclusion

- The model is very simple, but scored 0.999998 on the [Private Leaderboard](#).
- Such a high score is possible because the mass can easily be reconstructed and allows us to almost perfectly distinguish between signal and background. The evaluation metric used (weighted area under the ROC curve) helps too.
- At the same time, the correlation test can be passed by combining a strong (kinematic) prediction with a weaker (geometric) one.
- XGBoost is a powerful BDT library that can certainly be useful for HEP.
- This was the first time when real (background) data from the LHC were provided in a public ML challenge. Combining real and simulated data is really difficult.
- I encourage the Organizers to make the datasets publicly available and to provide the rest of the kinematic features (azimuthal angles of p_0 , p_1 , p_2).
- We all look forward to new Kaggle challenges from CERN and Yandex!

Thank you for your attention!