# Machine learning approach to analyze cosmics data

**Lia Lavezzi**
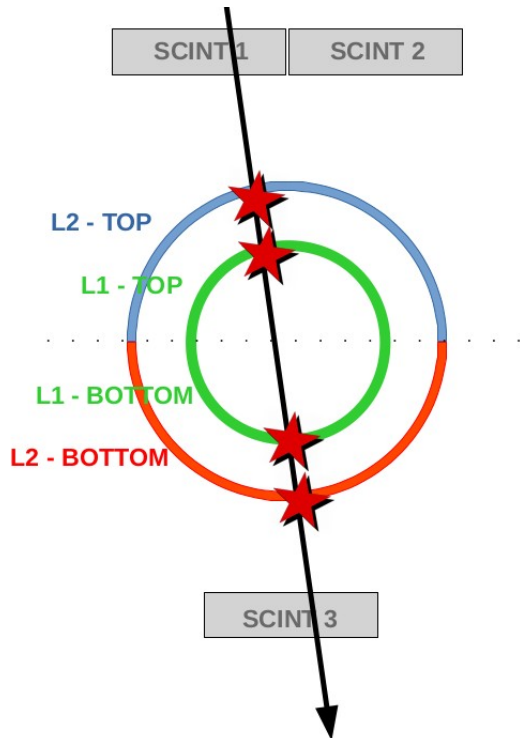
Università degli Studi di Torino / INFN–TO
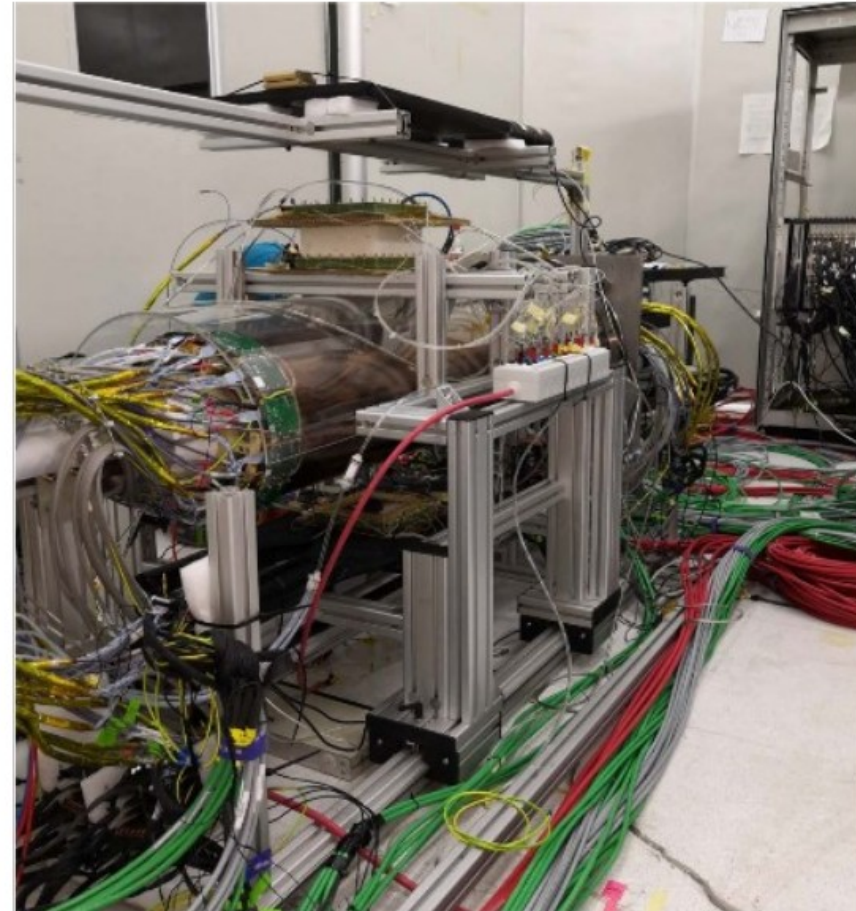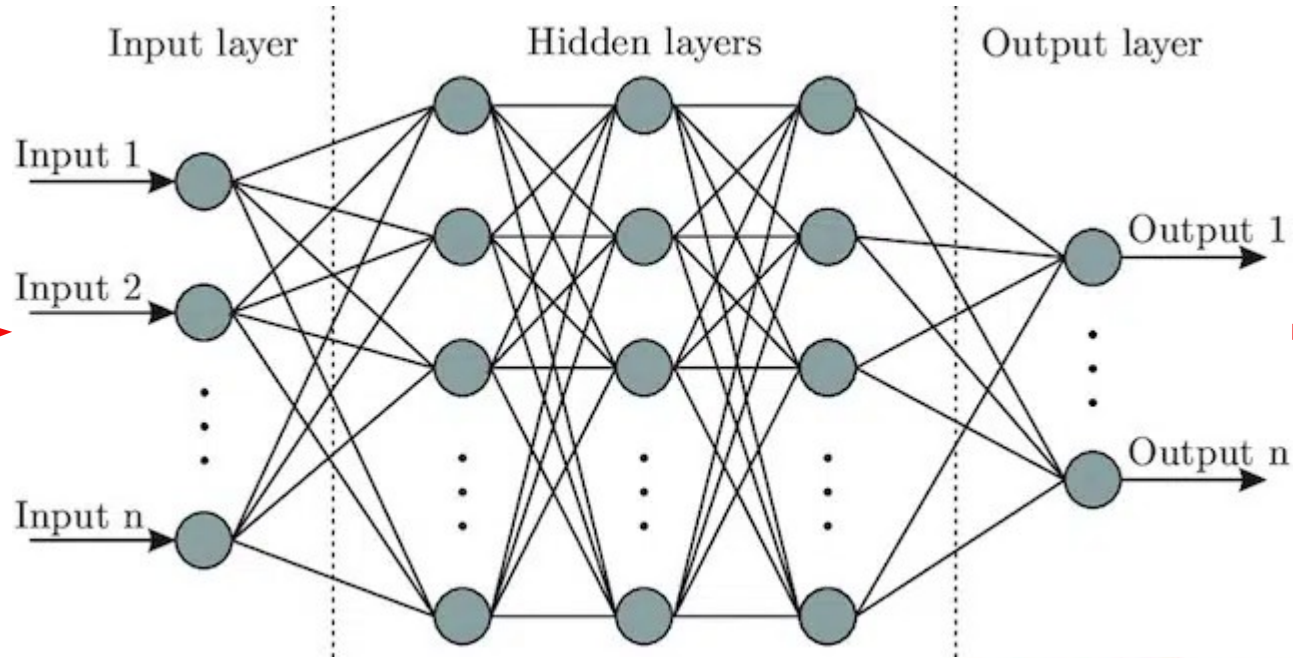
- Machine Learning to separate signal from noise at hit/cluster level
- Ongoing thesis @ Ferrara



SETUP - LAYER1 + LAYER2

Machine Learning

Measured quantities:
- hit/cluster charge
- signal start time, length ...
- charge sharing between x/v
- ...

Classification:
- Signal
- Noise

# Signal and Noise

The network must be **trained** to **learn** from **true** and **known** events

- **Signal:** golden run, no. 17
  Event selection with:
  - cgemboss.6.6.5.g
  - alignment: fixed
  - 4 plane fit → all 4 planes firing
  - no more than 100 cluster–2D
  - no chi2 cut
  - 129357 events

- **Noise:** run 44–49
  - hit in trigger time window
  - 98363 events

| RUN ID | | | | |
|---|---|---|---|---|
| **AQUIRED** | **SHARED** | **GAIN** | **THRSH** | **# EVENTS** |
| 351 | 17 | std | std | 234154 |
| 421 | 44 | off (HV=0) | std | 18853 |
| 422 | 45 | off (HV=0) | std | 3021 |
| 423 | 46 | off (HV=0) | std | 21168 |
| 424 | 47 | off (HV=0) | std | 24288 |
| 426 | 48 | off (HV=0) | std | 14469 |
| 428 | 49 | off (HV=0) | std | 16564 |

Run **44 − 49**: hit charge *vs* hit time shows a **structure** close to the event window

Without run **45 − 46** the structure disappears → exclude them
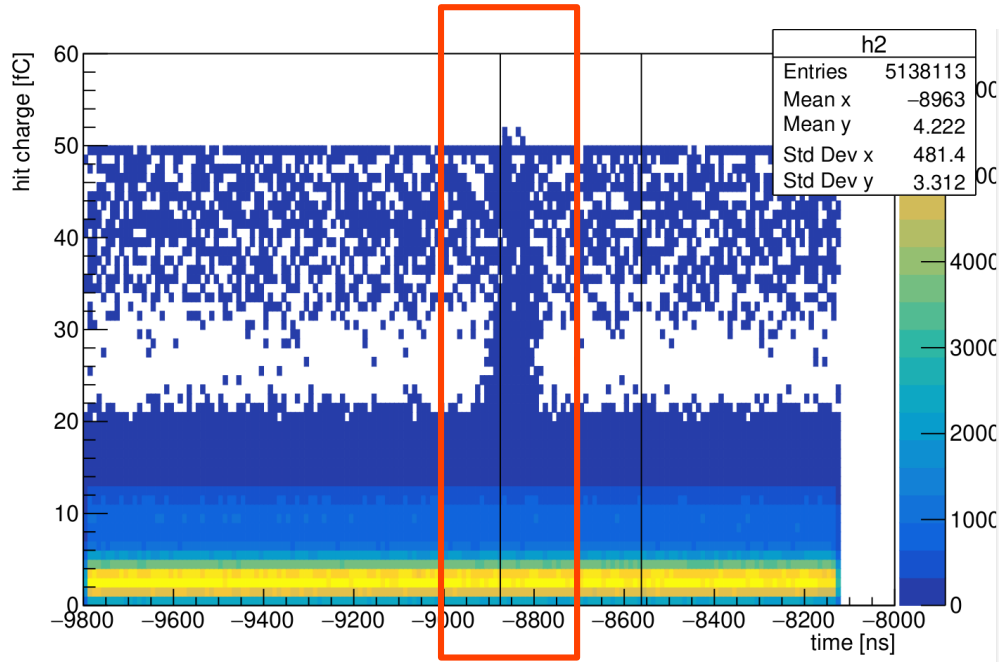
# NOISE CLEANING

The network must be **trained** to **learn** from **true** and **known** events

- **Signal:** golden run, no. 17
  Event selection with:
  - cgemboss.6.6.5.g
  - alignment: fixed
  - 4 plane fit → all 4 planes firing
  - no more than 100 cluster–2D
  - no chi2 cut
  - 129357 events

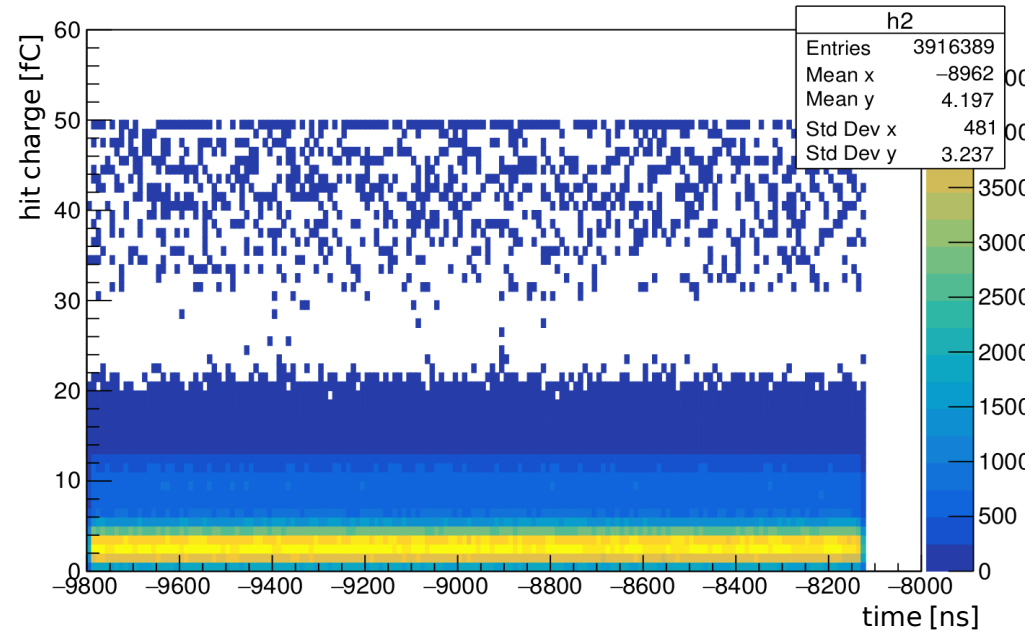- **Noise:** run 44, 47, 48, 49
  - hit in trigger time window
  - 74174 events

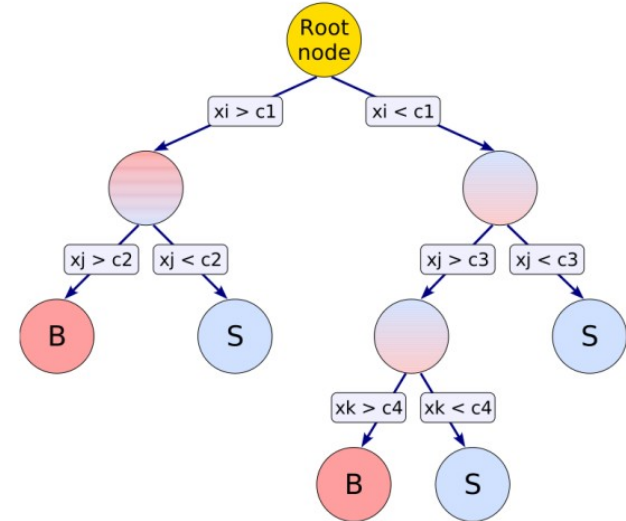| RUN ID | | | | |
|---|---|---|---|---|
| **AQUIRED** | **SHARED** | **GAIN** | **THRSH** | **# EVENTS** |
| 351 | 17 | std | std | 234154 |
| 421 | 44 | off (HV=0) | std | 18853 |
| 422 | 45 | off (HV=0) | std | 3021 |
| 423 | 46 | off (HV=0) | std | 21168 |
| 424 | 47 | off (HV=0) | std | 24288 |
| 426 | 48 | off (HV=0) | std | 14469 |
| 428 | 49 | off (HV=0) | std | 16564 |

- Toolkit for Multi Variate Analysis
- C++ Object Oriented / ROOT
- Many algorithms available:

|  | | Cuts | Likeli-hood | PDE-RS / k-NN | PDE-Foam | H-Matrix | Fisher / LD | MLP | BDT | Rule-Fit | SVM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | |
| Perfor-mance | No or linear correlations | ★ | ★★ | ★ | ★ | ★ | ★★ | ★★ | ★ | ★★ | ★ |
| | Nonlinear correlations | ○ | ○ | ★★ | ★★ | ○ | ○ | ★★ | ★★ | ★★ | ★★ |
| Speed | Training | ○ | ★★ | ★★ | ★★ | ★★ | ★★ | ★ | ★ | ★ | ○ |
| | Response | ★★ | ★★ | ○ | ★ | ★★ | ★★ | ★★ | ★ | ★★ | ★ |
| Robust-ness | Overtraining | ★★ | ★ | ★ | ★ | ★★ | ★★ | ★ | ★ | ★ | ★★ |
| | Weak variables | ★★ | ★ | ○ | ○ | ★★ | ★★ | ★ | ★★ | ★ | ★ |
| Curse of dimensionality | | ○ | ★★ | ○ | ○ | ★★ | ★★ | ★ | ★ | ★ | |
| Transparency | | ★★ | ★★ | ★ | ★ | ★★ | ★★ | ○ | ○ | ○ | ○ |

The column header "MVA METHOD" spans the method columns.

**Table 6**: Assessment of MVA method properties. The symbols stand for the attributes "good" (★★), "fair" (★) and "bad" (○). "Curse of dimensionality" refers to the "burden" of required increase in training statistics and processing time when adding more input variables. See also comments in the text. The FDA method is not listed here since its properties depend on the chosen function.

*Most straightforward one for starters*

### BOOSTED DECISION TREE



- Series on nodes: @ each node a cut decides *left/right*
- End nodes (*leafs*) are signal/background
- Several trees with different *weights* to form a *forest*
- BDT outcome is the response of the forest
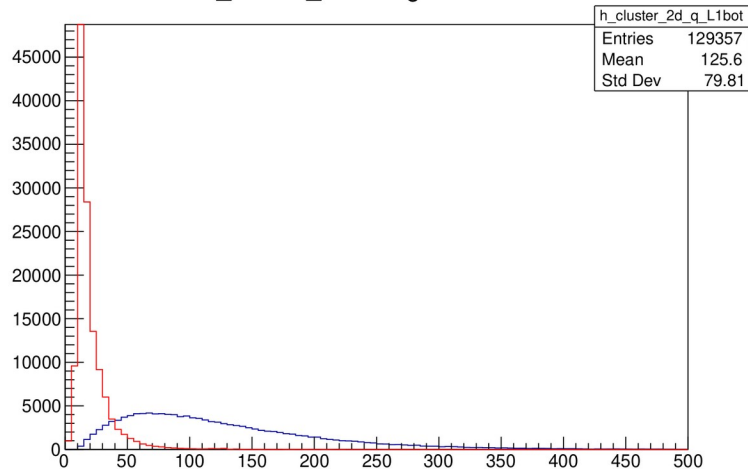
$$\hat{y}(x) = \sum_t w_t h_t(x)$$

# L1 BOTTOM

**SIGNAL**

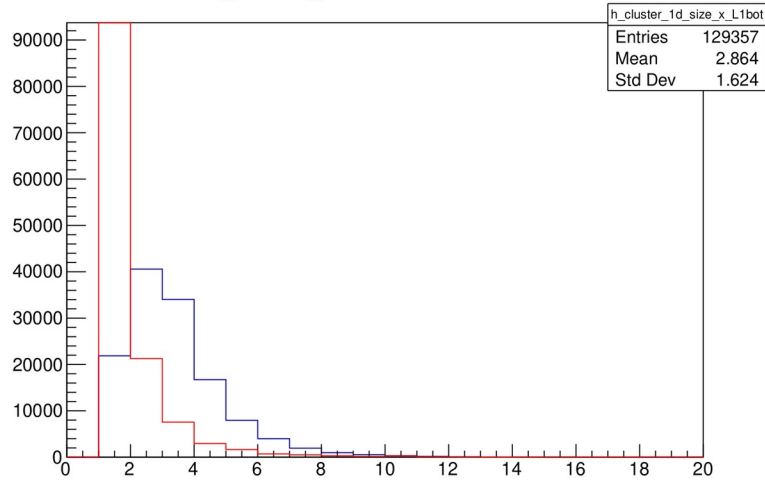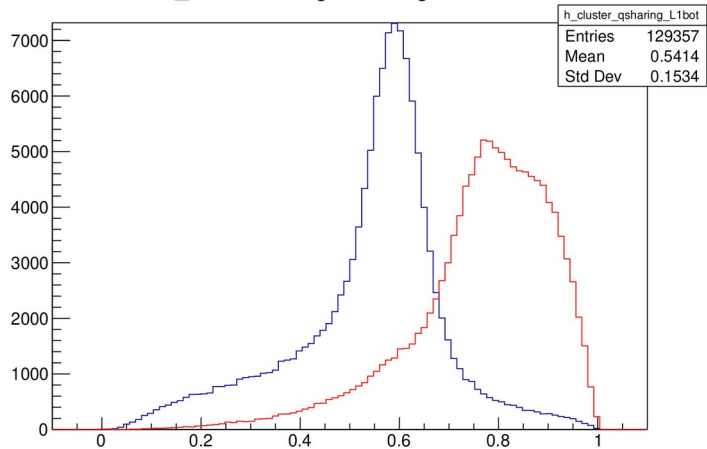**NOISE**



h_cluster_2d charge on L1 bot

| h_cluster_2d_q_L1bot | |
|---|---|
| Entries | 129357 |
| Mean | 125.6 |
| Std Dev | 79.81 |



h_cluster_1d cl. size on L1 bot x

| h_cluster_1d_size_x_L1bot | |
|---|---|
| Entries | 129357 |
| Mean | 2.864 |
| Std Dev | 1.624 |



h_cluster charge sharing on L1 bot

| h_cluster_qsharing_L1bot | |
|---|---|
| Entries | 129357 |
| Mean | 0.5414 |
| Std Dev | 0.1534 |



h_hit_fast_t on L1 bot x

| h_hit_fast_t_L1bot_x | |
|---|---|
| Entries | 129357 |
| Mean | −8740 |
| Std Dev | 39.28 |



h_hit_delta_t on L1 bot v

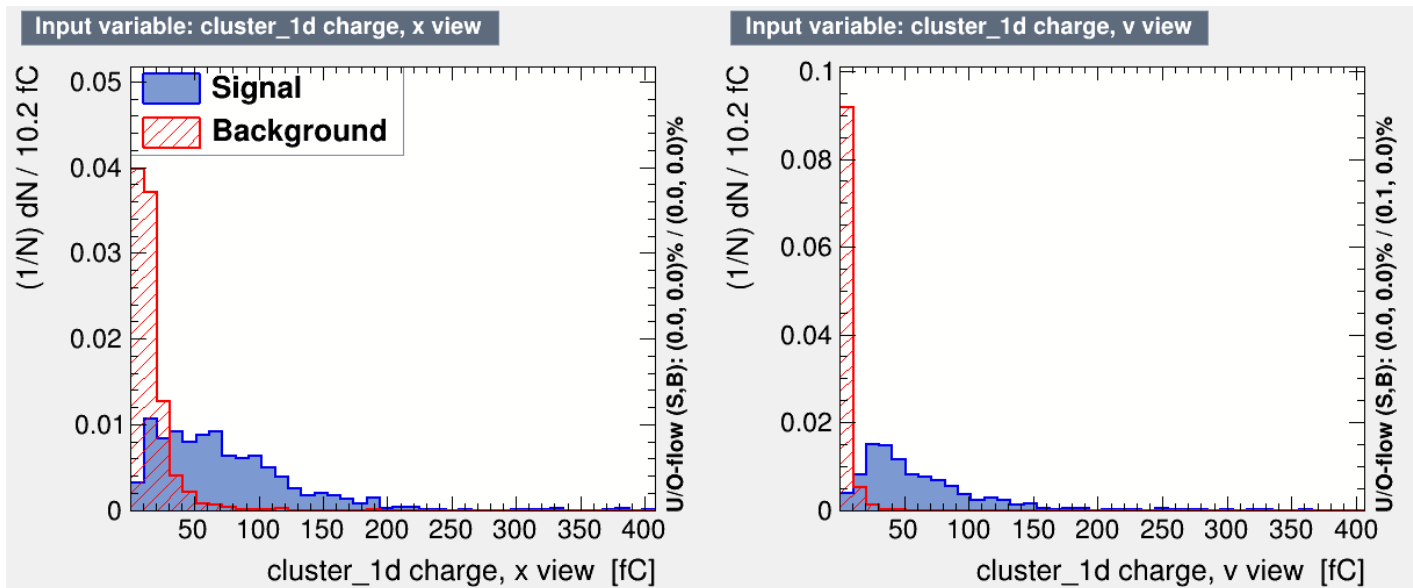| h_hit_delta_t_L1bot_v | |
|---|---|
| Entries | 370768 |
| Mean | 38.87 |
| Std Dev | 37.43 |

**Very first** attempt:

```
dataloader->AddVariable( "mva_cluster_1d_qx", "cluster_1d charge, x view", "fC", 'F');
dataloader->AddVariable( "mva_cluster_1d_qv", "cluster_1d charge, v view", "fC", 'F');
```

- #signal events for training = 1000
- #background events for training = 1000
- split mode = random
- method: BDT
- NTrees=850
- MinNodeSize=2.5%
- MaxDepth=3
- BoostType=AdaBoost
- AdaBoostBeta=0.5
- UseBaggedBoost
- BaggedSampleFraction=0.5
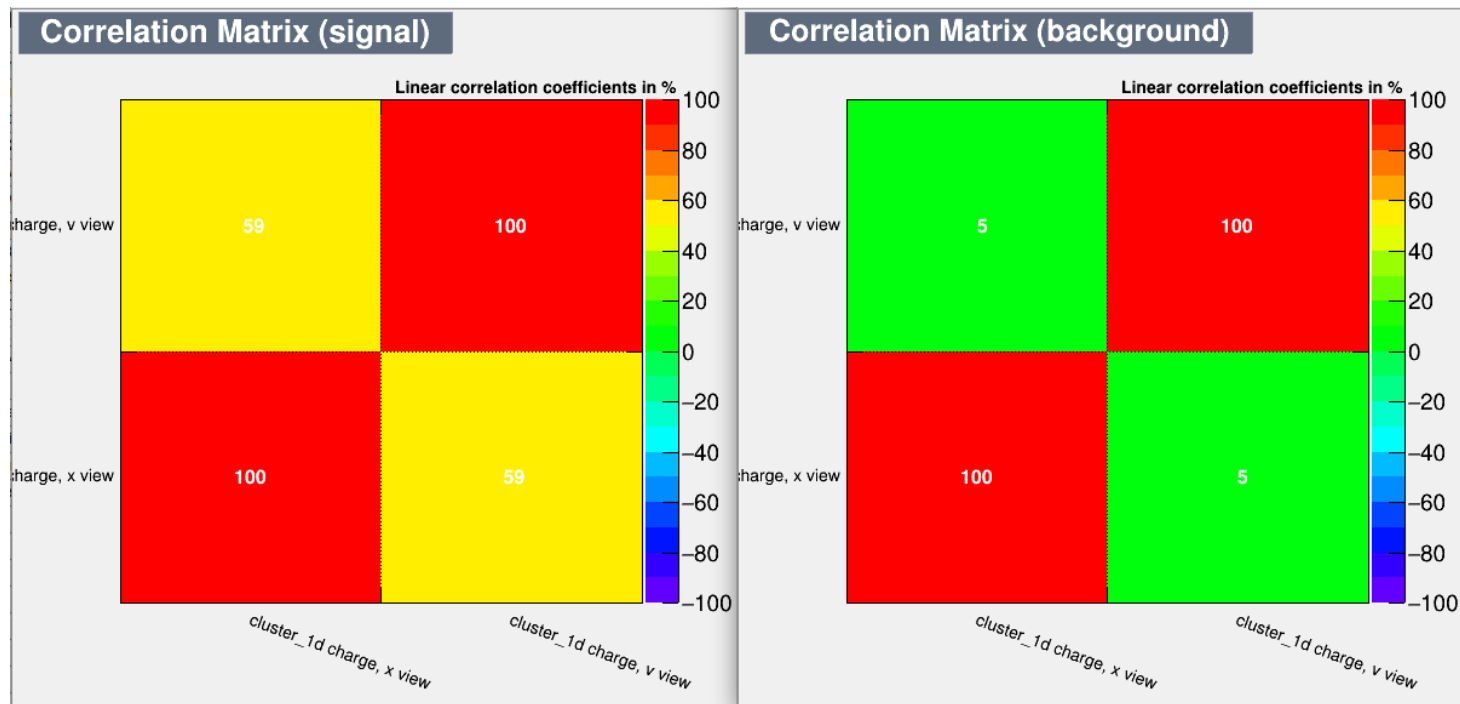- SeparationType=GiniIndex
- nCuts=20"

### VARIABLE VALUES

**Very first** attempt:

```
dataloader->AddVariable( "mva_cluster_1d_qx", "cluster_1d charge, x view", "fC", 'F');
dataloader->AddVariable( "mva_cluster_1d_qv", "cluster_1d charge, v view", "fC", 'F');
```

- #signal events for training = 1000
- #background events for training = 1000
- split mode = random
- method: BDT
- NTrees=850
- MinNodeSize=2.5%
- MaxDepth=3
- BoostType=AdaBoost
- AdaBoostBeta=0.5
- UseBaggedBoost
- BaggedSampleFraction=0.5
- SeparationType=GiniIndex
- nCuts=20"

## CORRELATIONS

**Very first** attempt:

```
dataloader->AddVariable( "mva_cluster_1d_qx", "cluster_1d charge, x view", "fC", 'F');
dataloader->AddVariable( "mva_cluster_1d_qv", "cluster_1d charge, v view", "fC", 'F');
```

- #signal events for training = 1000
- #background events for training = 1000
- split mode = random
- method: BDT
- NTrees=850
- MinNodeSize=2.5%
- MaxDepth=3
- BoostType=AdaBoost
- AdaBoostBeta=0.5
- UseBaggedBoost
- BaggedSampleFraction=0.5
- SeparationType=GiniIndex
- nCuts=20"

SELECT CUT

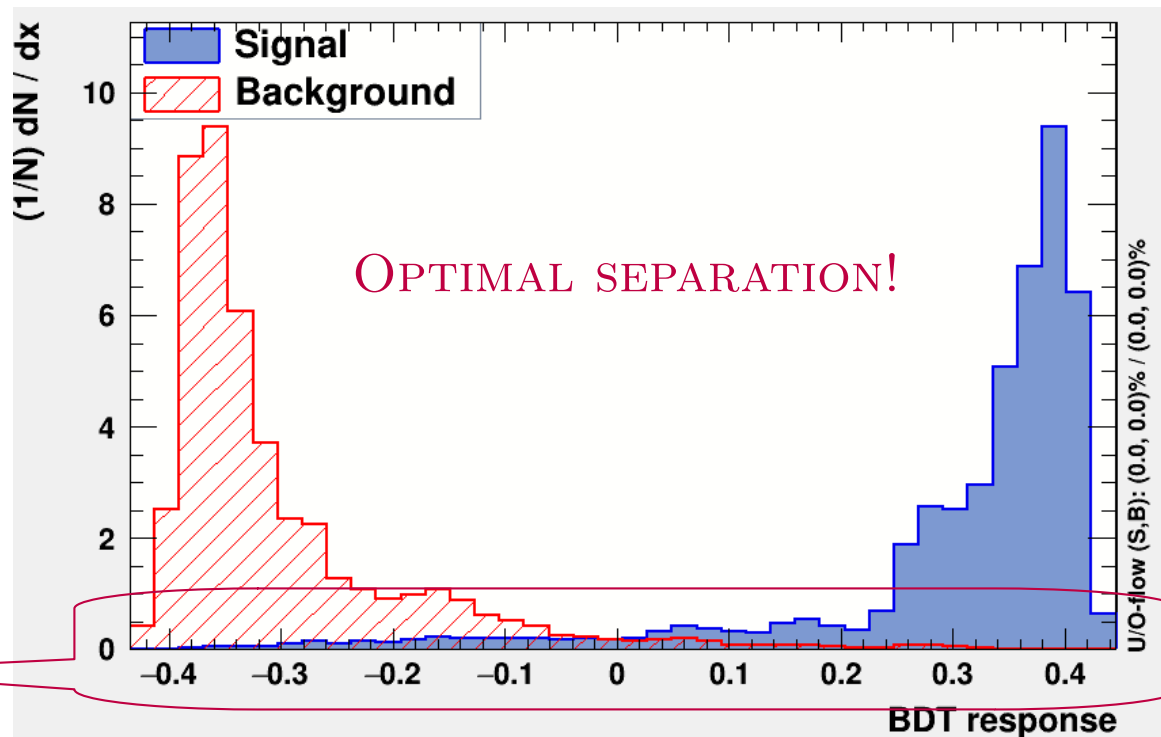OPTIMAL SEPARATION!



Machine Learning

**Very first** attempt:

```
dataloader->AddVariable( "mva_cluster_1d_qx", "cluster_1d charge, x view", "fC", 'F');
dataloader->AddVariable( "mva_cluster_1d_qv", "cluster_1d charge, v view", "fC", 'F');
```

- #signal events for training = 1000
- #background events for training = 1000
- split mode = random
- method: BDT
- NTrees=850
- MinNodeSize=2.5%
- MaxDepth=3
- BoostType=AdaBoost
- AdaBoostBeta=0.5
- UseBaggedBoost
- BaggedSampleFraction=0.5
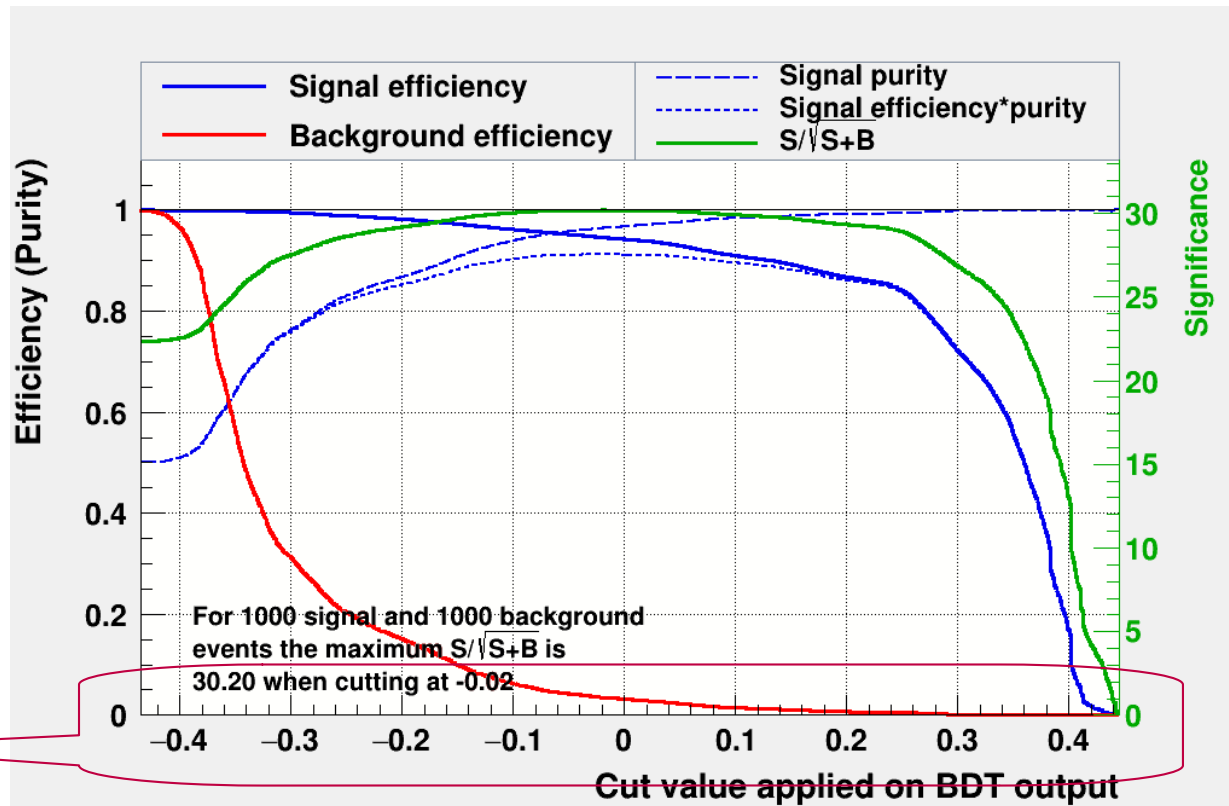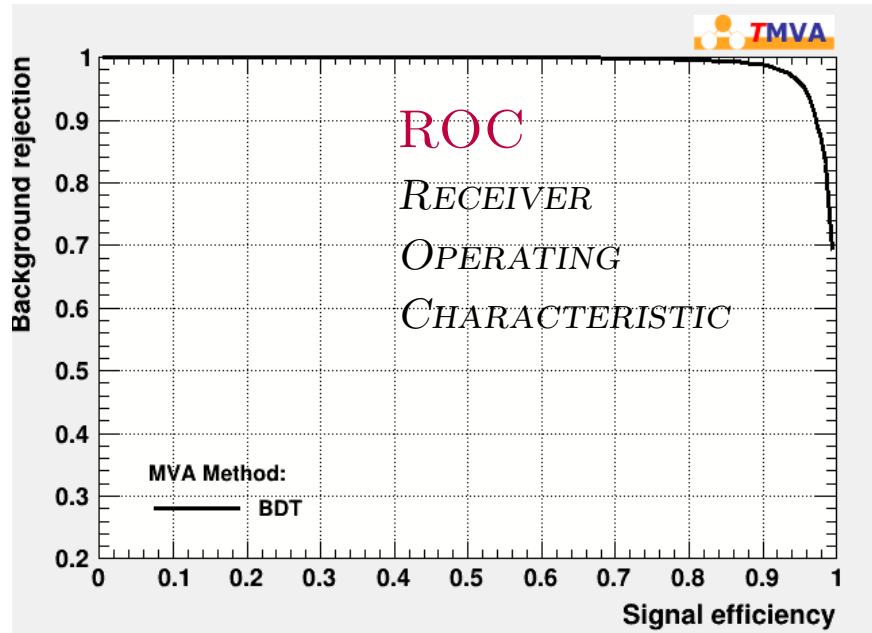- SeparationType=GiniIndex
- nCuts=20"



SELECT CUT

Machine Learning

**ROC**

*RECEIVER*

*OPERATING*

*CHARACTERISTIC*

MVA Method:

BDT

We want **good noise rejection** and **good signal efficiency** → give–and–take

- efficiency = 1 − rejection
- given @ three background eff. Levels

example:

@B=0.10 → background rejection = 0.90 → signal efficiency = 0.972

```
Testing efficiency compared to training efficiency (overtraining check)
-----------------------------------------------------------------------------------
DataSet           MVA             Signal efficiency: from test sample (from training sample)
Name:             Method:         @B=0.01               @B=0.10              @B=0.30
-----------------------------------------------------------------------------------
dataset           BDT             : 0.888 (0.902)       0.972 (0.974)        0.994 (0.994)
-----------------------------------------------------------------------------------
```

# SUMMARY

- The whole procedure works
- Different functions have been understood
- Signal and noise training/testing samples have been updated

## TO DO LIST

- Optimize the BDT parameters
- Test other variables
- Test other ML methods

# SUMMARY

- The whole procedure works
- Different functions have been understood
- Signal and noise training/testing samples have been updated

## TO DO LIST

- Optimize the BDT parameters
- Test other variables
- Test other ML methods