# Machine Learning For Particle Identification

05.June.2018 I  Waleed Esmail

Institut für Kernphysik (IKP), Forschungszentrum Jülich
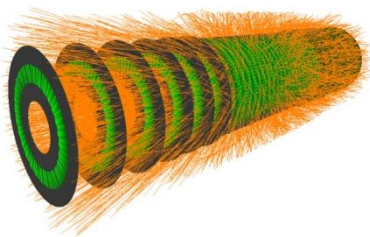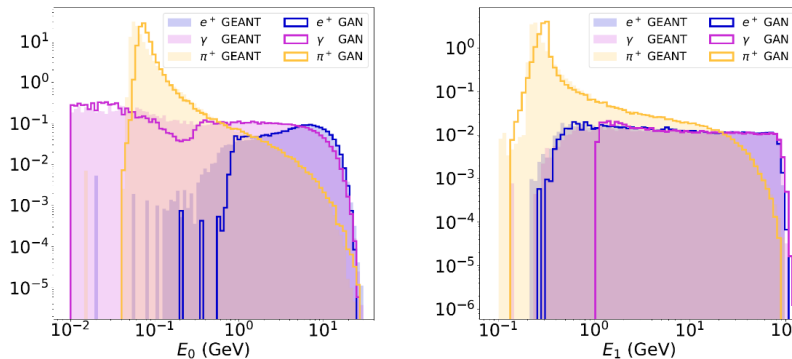
# Outlines:

➤ Whetting Your Appetite.

➤ Introduction to Machine Learning.

➤ Motivating the Key Concepts.

➤ Boosted Decision Trees.

➤ Artificial Neural Networks.

➤ Integration to PandaRoot.

➤ Conclusion.

JÜLICH
Forschungszentrum

# Whetting your Appetite:

- Why Multivariate Analysis (Machine Learning Techniques)?

  - Classification (Higgs discovery).

  - Clustering (Tracking).

  - Pattern Recognition (Tracking and Jet Images).
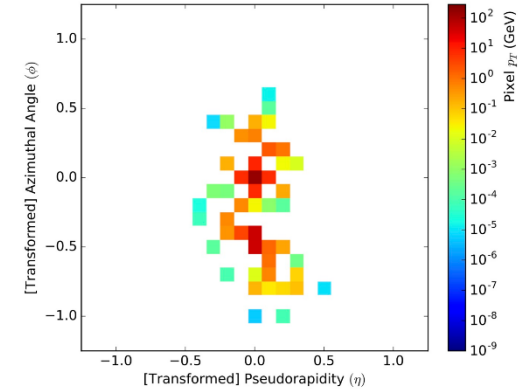
  - Generative Learning GAN (Simulation "GeantV").

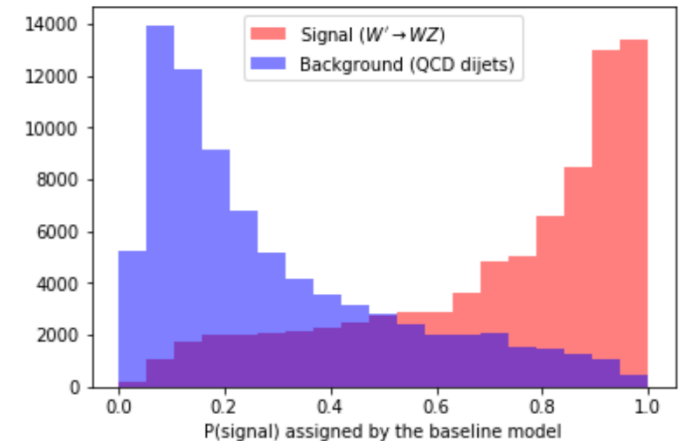**speed-up factors of up to 100,000x**



Jet Image using CNN.



Total prizes:
$25,000

Next important deadline:
August 6, 2018 - Entry deadline. You must accept the competition rules before this date in order to compete.

Join the competition

JÜLICH
Forschungszentrum

# Introduction to Machine Learning:

- Machine Learning (ML) is about modeling your data.

- Developing *self-learning* algorithms to gain knowledge from data in order to make *predictions*.

- There are *three* types of ML *algorithms*:

  1. **Supervised Learning:**

     - Learn a model from labeled training data.

     - Classification (discrete labels), and Regression (continuous response variable).
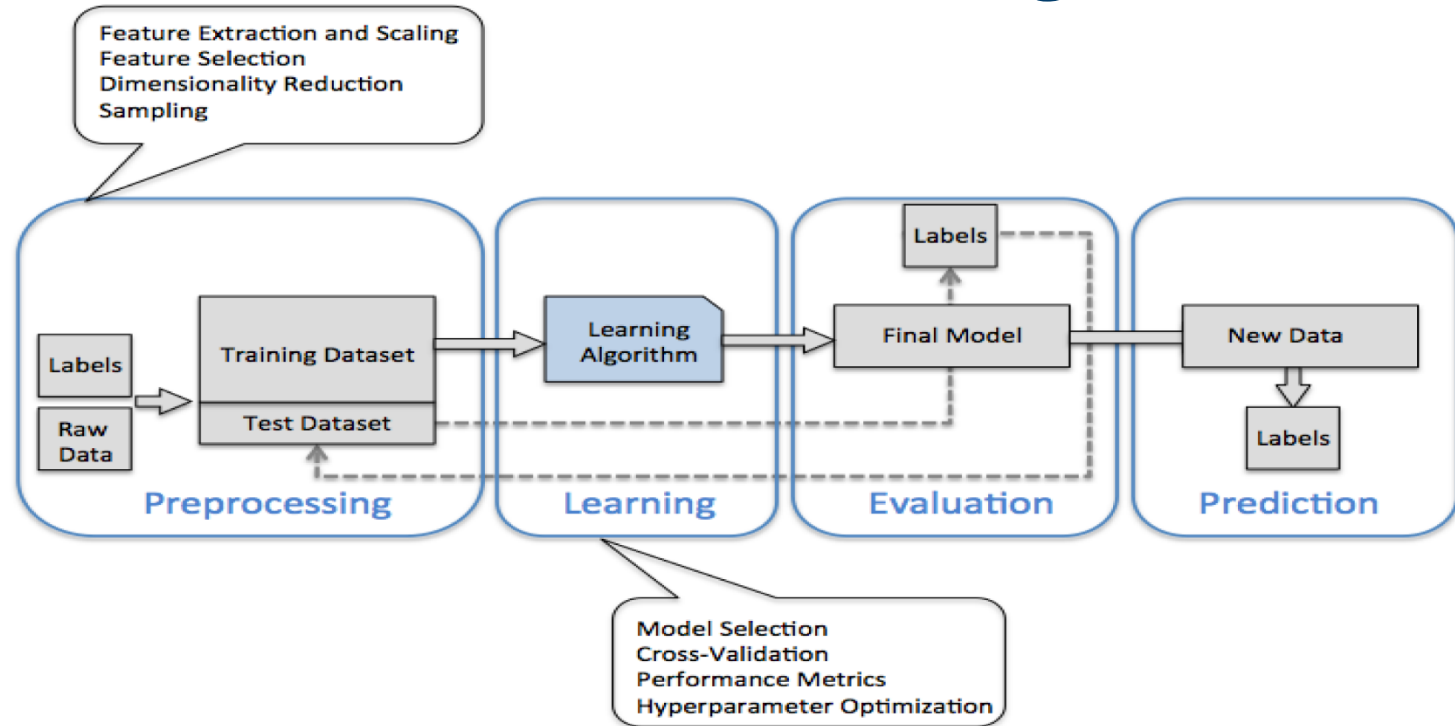
  2. **Unsupervised Learning:**

     - We deal with unlabeled data to explore its structure.

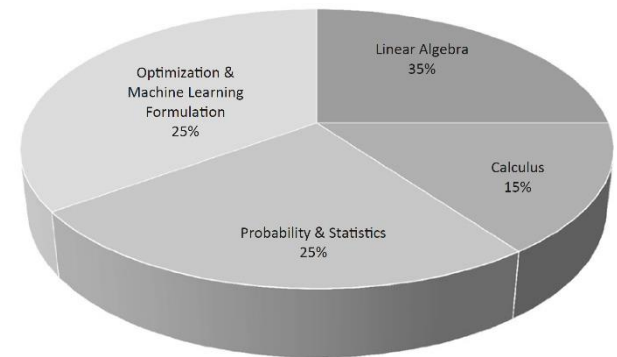     - Clustering (subgrouping), and Dimensionality Reduction (feature preprocessing.)

  3. **Reinforcement learning:**

     - Develop a system (agent) to improve performance via reward maximization.

**JÜLICH**
Forschungszentrum

# Introduction to Machine Learning:



Strong grasp of the core concepts of mathematics enables one to ***select the right algorithm***, Also, it enables one to ***tune*** machine-learning/deep-learning models better.
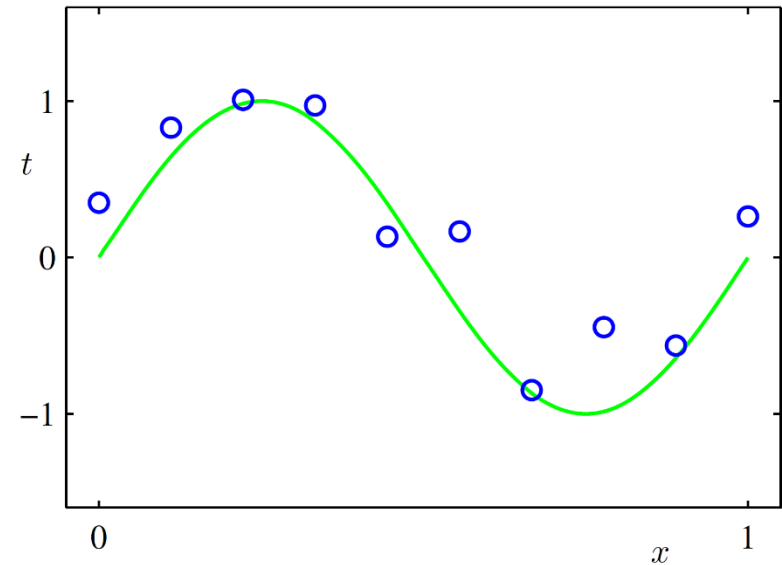
# Motivating the Key Concepts:

- A simple *regression* problem.

- We observe a real-valued input variable ($x$), and we wish to *predict* the value of a real-valued target variable ($t$).

- The *green* curve shows the function $\sin(2\pi x)$ used to generate the data.

- The *blue* points are obtained by adding Gaussian noise.

- Fit the data using a *polynomial* function:

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$
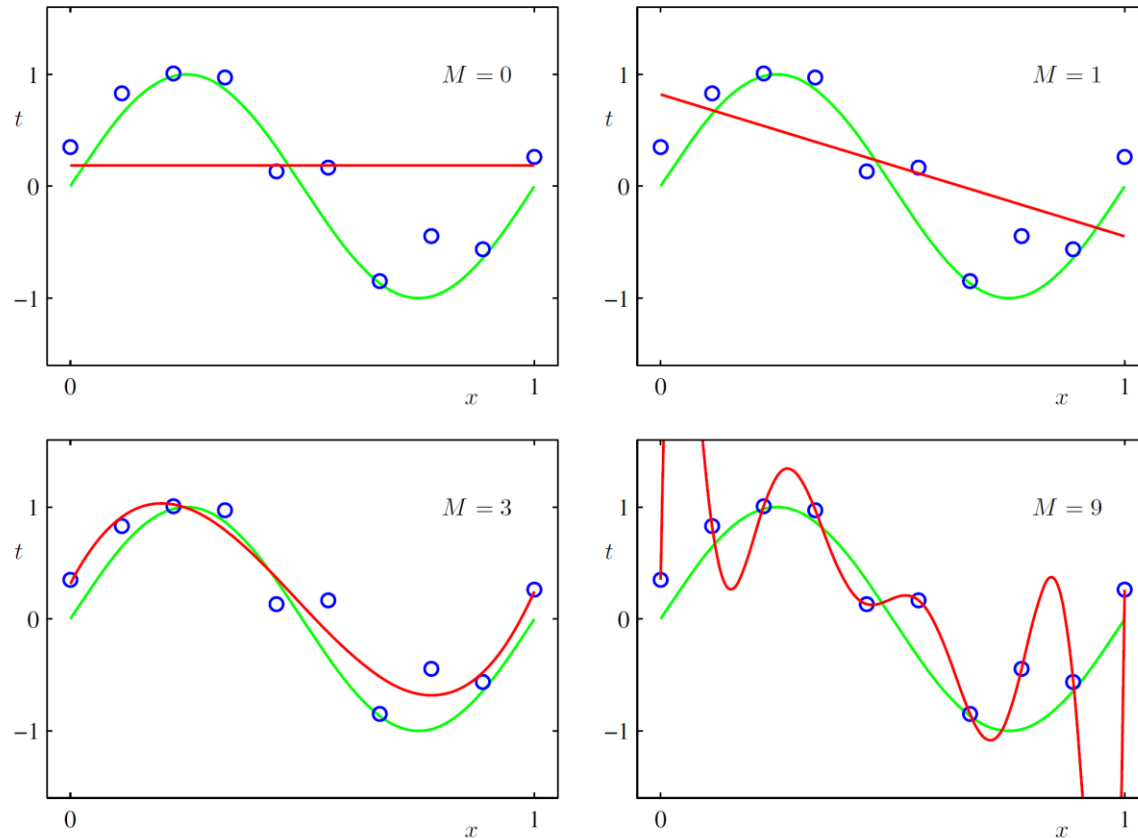
- Polynomial degree M ➡ (*model selection*).

- *Weights* can be determined by fitting polynomial to the training data, by *minimizing* an *error function.*

**JÜLICH**
Forschungszentrum

# Motivating the Key Concepts:

- One common choice is the **_root-mean-square_**.
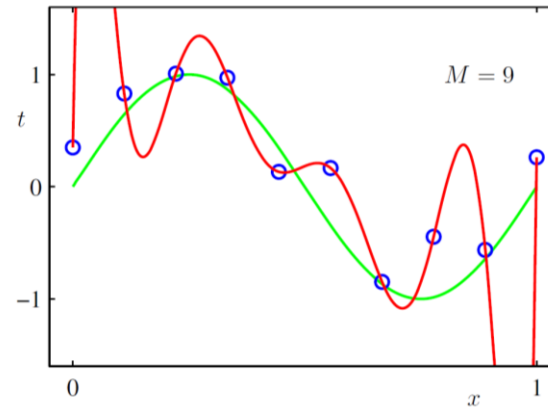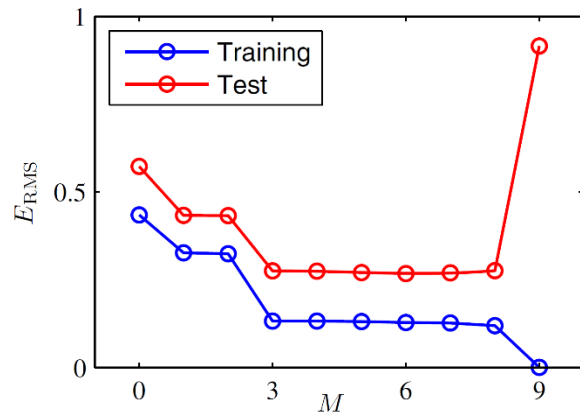
$$E_{\mathrm{RMS}} = \sqrt{2E(\mathbf{w})/N} \qquad \Longrightarrow \qquad E(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}\{y(x_n, \mathbf{w}) - t_n\}^2$$

**JÜLICH**
Forschungszentrum

# Motivating the Key Concepts:

- One common choice is the **root-mean-square**.

$$E_{\mathrm{RMS}} = \sqrt{2E(\mathbf{w})/N}$$

$$\Longrightarrow \qquad E(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}\{y(x_n, \mathbf{w}) - t_n\}^2$$



**Over-fitting**

- But how the weights of the model can be computed. **Probability** theory comes to rescue.

JÜLICH
Forschungszentrum

# Motivating the Key Concepts:

- Given $(x)$, assume that the target variable $(t)$ has a Gaussian distribution:

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}\left(t|y(x, \mathbf{w}), \beta^{-1}\right) \qquad \beta^{-1} = \sigma^2$$

- Determine the unknown parameters by ***maximum likelihood method.***

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}\left(t_n|y(x_n, \mathbf{w}), \beta^{-1}\right)$$

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2}\sum_{n=1}^{N}\{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi)$$
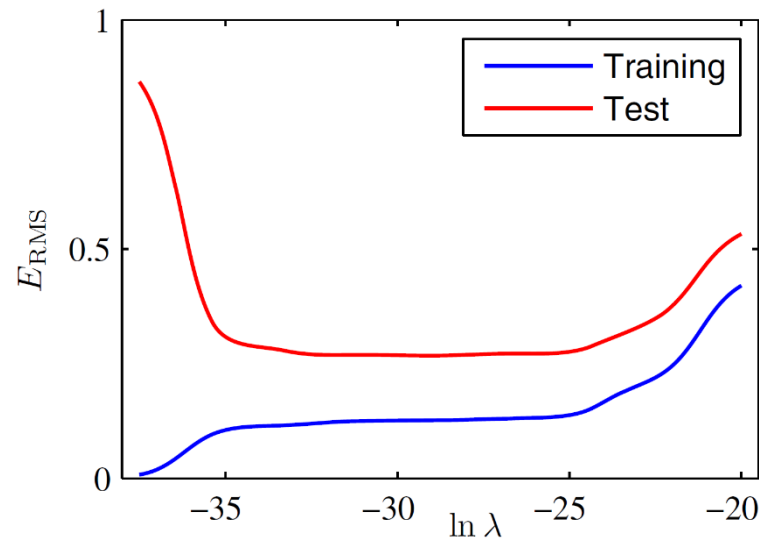
- By maximizing with respect to $(w)$, we obtain the desired solution $(w_{ML})$.

- Finally make new predictions.

JÜLICH
Forschungszentrum

# Motivating the Key Concepts:

- Back to over-fitting issue, one technique to control it is called **regularization**.

- Regularization involves adding **penalty term** to the error function.
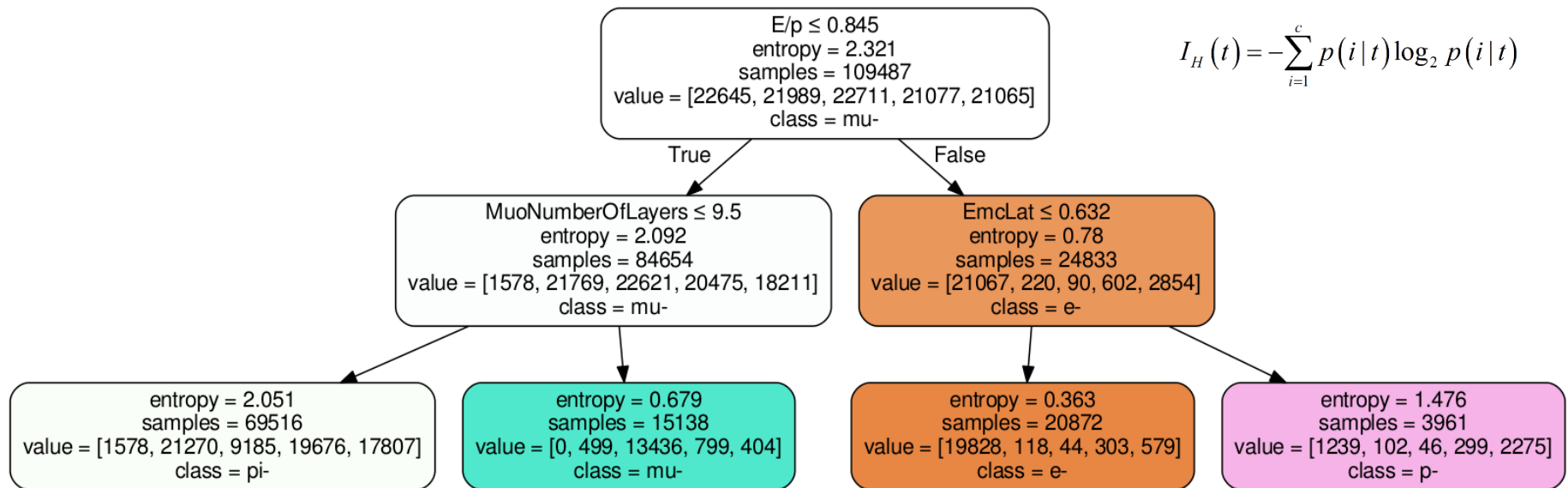
$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- $\lambda$ governs the relative importance of the regularization term compared to the error term.

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):

- Idea: ***break down the data*** by making decisions based on the features in the training set.
- The splitting is done by maximizing the **Information Gain** $IG(D_p, f) = I(D_p) - \sum_{j=1}^{m} \frac{N_j}{N_p} I(D_j)$
- **BDT** is an **ensemble method.** The key concept behind **boosting** is to focus on training samples that are hard to classify.

$$I_H(t) = -\sum_{i=1}^{c} p(i|t) \log_2 p(i|t)$$

```
                          E/p ≤ 0.845
                          entropy = 2.321
                          samples = 109487
            value = [22645, 21989, 22711, 21077, 21065]
                          class = mu-
```

True / False

```
   MuoNumberOfLayers ≤ 9.5                  EmcLat ≤ 0.632
      entropy = 2.092                        entropy = 0.78
      samples = 84654                        samples = 24833
value = [1578, 21769, 22621, 20475, 18211]  value = [21067, 220, 90, 602, 2854]
        class = mu-                              class = e-
```

```
   entropy = 2.051          entropy = 0.679          entropy = 0.363          entropy = 1.476
   samples = 69516          samples = 15138          samples = 20872          samples = 3961
value = [1578, 21270,    value = [0, 499, 13436,  value = [19828, 118,     value = [1239, 102,
9185, 19676, 17807]       799, 404]                44, 303, 579]            46, 299, 2275]
   class = pi-              class = mu-              class = e-               class = p-
```

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):

- Two event generators are used for training.

## 1. BoxGenerator:

- **momentum range**: (0.2 - 5 ) GeV.

- **phi range**: 0 - 360$^o$.

- **theta range**: 0 - 180$^o$.

- **particle species**: [ e$^{\mp}$, $\pi^{\mp}$, $\mu^{\mp}$, k$^{\mp}$, p$^{\mp}$ ]. One particle per event.

## 2. EvtGen:

- $p\bar{p} \rightarrow X\bar{X}Y\bar{Y}$, where X, and Y = e$^{\mp}$, $\pi^{\mp}$, $\mu^{\mp}$, k$^{\mp}$, p$^{\mp}$

- **Beam momentum:** 15 GeV/c.

*Particles are matched to their **MC truth** information.*

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):

- Input Features:



Feature Correlations (EvtGen)

Feature Correlations (Box)

# Boosted Decision Trees (BDT):

- Data are organized into python **DataFrame**. Data was splitted into **training** (**70%**) and **testing** (**30%**) sets.

```
+--------+--------+------+----------+-------+-------+-------+-----------+--------------+-----------+
| energy|momentum|charge|  position|MvdHits|GemHits|SttHits|TofStopTime|TofTrackLength|EmcCalEnergy|
+--------+--------+------+----------+-------+-------+-------+-----------+--------------+-----------+
|1.45992|  2.1119|    -1|1.53556E-4|      3|      0|     26|        0.0|           0.0|   0.250083|
|4.14557| 17.1663|    -1|6.65813E-6|      5|      0|     17|    3.95565|       117.048|   0.991413|
|3.51102| 12.3078|    -1|0.00648225|      1|      0|     24|        0.0|           0.0|    2.21215|
|3.45948| 11.9486|    -1| 5.4671E-6|      4|      0|     22|    2.89786|       86.4262|    0.29421|
|4.78585| 22.8849|    -1|6.36045E-5|      3|      0|     26|        0.0|           0.0|    2.62603|
+--------+--------+------+----------+-------+-------+-------+-----------+--------------+-----------+
```



**Trained on Evt**

**Trained on Box**

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):

## *1. Receiver Operating Characteristic (ROC).*

**Trained on Evt**

**Trained on Box**



True Positives (TP) & True Negatives (TN).
False Positives (FP) & False Negatives (FN).

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):

## *2. Precision, and Recall.*

**Trained on Evt**

**Trained on Box**



$$\text{recall} = \frac{TP}{TP+FN} \quad \text{precision} = \frac{TP}{TP+FP}$$

JÜLICH
Forschungszentrum

# Boosted Decision Trees (BDT):
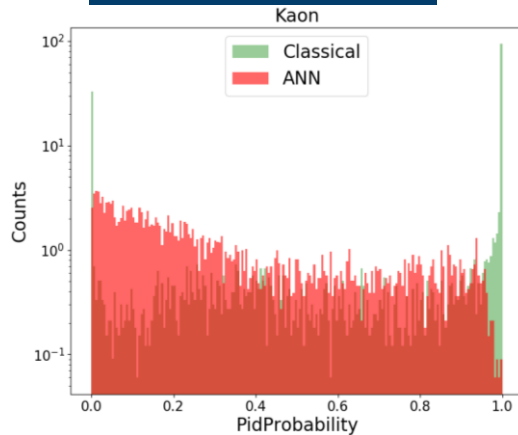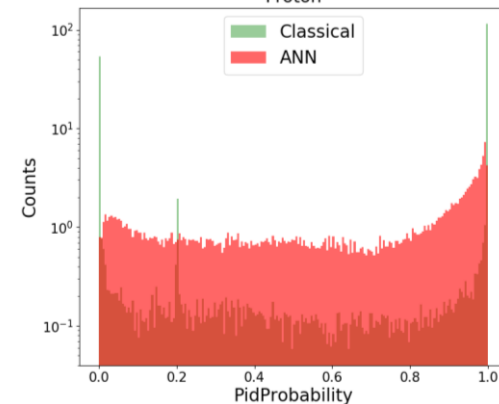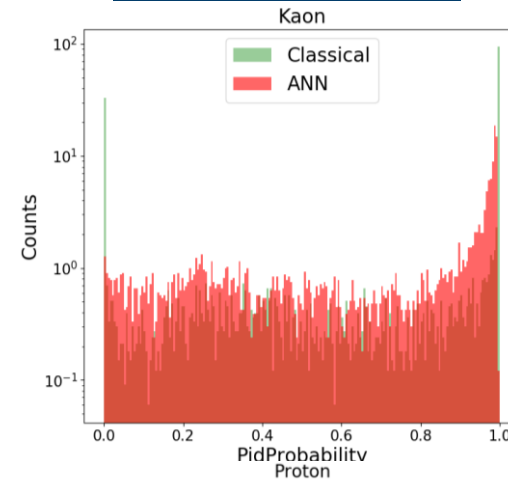
## *Particle Identification (PID) Probabilities.*

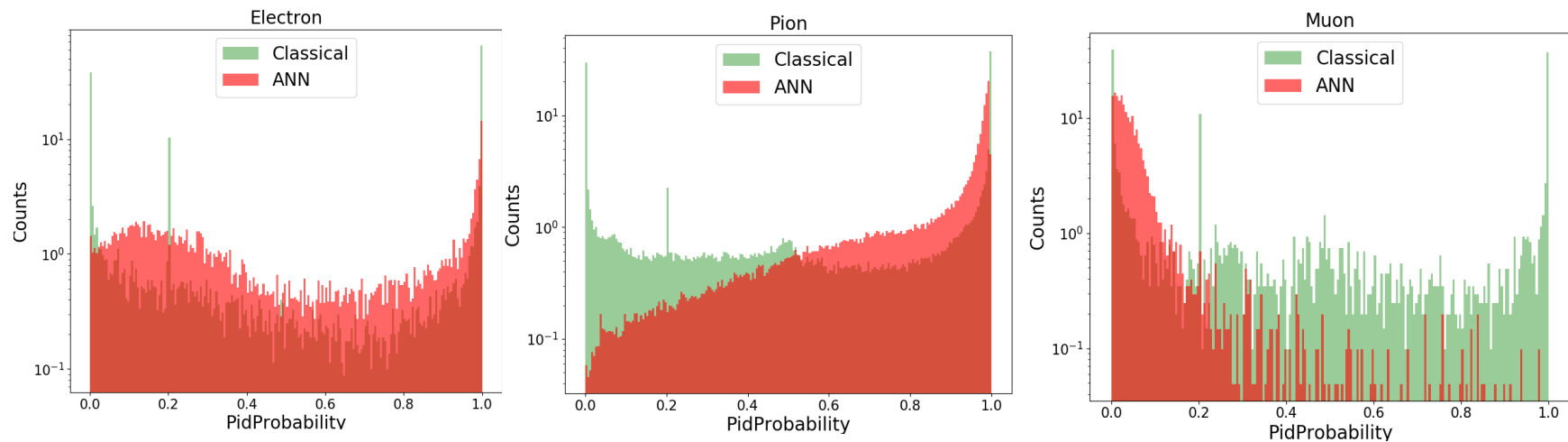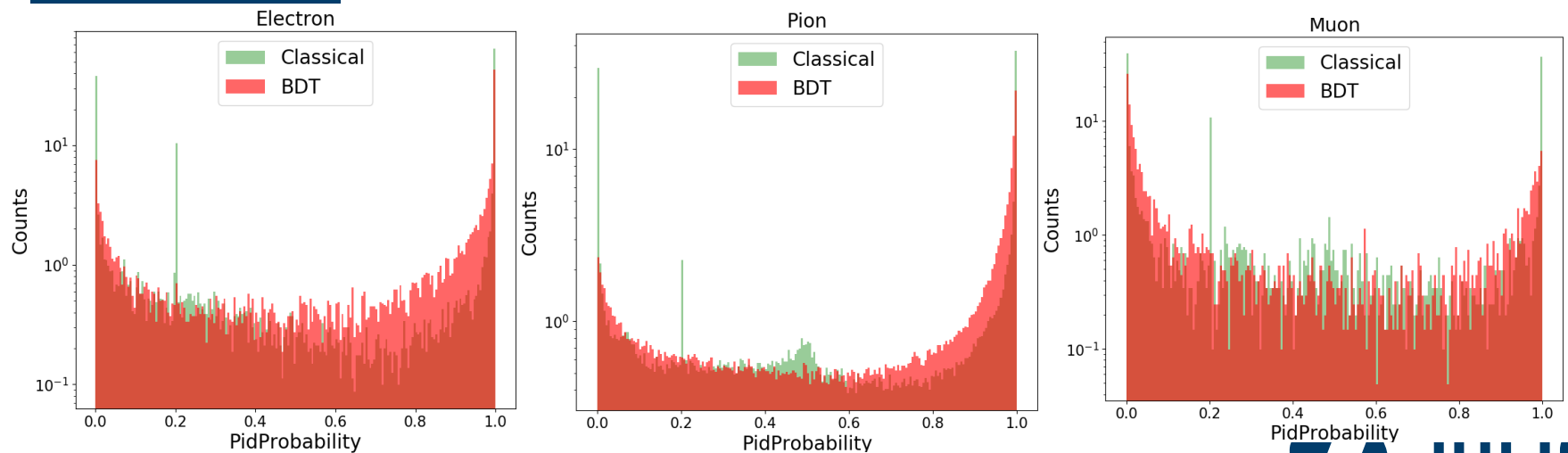We tested the trained algorithm on data generated by **DPM generator** (elastic + inelastic).

JÜLICH
Forschungszentrum

# Deep Learning:

## Artificial Neural Networks (ANN):

- Neural Networks is about meaningfully *transform the data.*



- *Training loop of DNN:*
  - ➤ *Input data*, parameterize by *weights* (transform the data), predict and compute the *loss score*, and update (*epochs*).

# Artificial Neural Networks (ANN):



- The update step is guided by ***minimizing the loss function*** by calculating its ***gradient***.
- Model parameters (weights) are updated through back-propagation algorithm.

- ***Keras*** Python Package was used.
- Keras is built on top of ***tensor-flow***, and can run on GPUs.
- Deep net with ***6 hidden layers***, about ~ 50,000 trainable parameters.
- ***L2 regularization*** is used to avoid ***overfitting***.

JÜLICH
Forschungszentrum

# Artificial Neural Networks (ANN):

# Artificial Neural Networks (ANN):

# Artificial Neural Networks (ANN):

## *Particle Identification (PID) Probabilities.*

We tested the trained algorithm on data generated by **DPM generator** (elastic + inelastic).

# Integration to PandaRoot:

.

```python
In [1]: from sklearn.externals import joblib
        clf = joblib.load('savedRF.pkl')
```

```python
In [5]: entry_to_predict_on = s_df.iloc[:,0:-1]
        particles = clf.predict_proba(entry_to_predict_on)
        List = ['e-','pi-','mu-','k-','p-']
        particles = pd.DataFrame(particles, columns=List)
        particles['event'] = s_df.loc[:,'events']
```

**Associator Task** ➡

output ⬇

PndPidMlAssociatorTask

JÜLICH
Forschungszentrum

# Integration to PandaRoot:

```cpp
FairRunAna* fRun = new FairRunAna();
FairRuntimeDb* rtdb = fRun->GetRuntimeDb();
fRun->SetInputFile(inPidFile);
fRun->AddFriend("signal_ml.root");
```

```cpp
theAnalysis->FillList(muplus,  "MuonTightPlus",  "PidAlgoMl");
theAnalysis->FillList(muminus, "MuonTightMinus", "PidAlgoMl");
theAnalysis->FillList(piplus,  "PionLoosePlus",  "PidAlgoMl");
theAnalysis->FillList(piminus, "PionLooseMinus", "PidAlgoMl");
```

JÜLICH
Forschungszentrum

# Integration to PandaRoot:

- [https://github.com/wesmail/MLPID_For_PANDARoot](https://github.com/wesmail/MLPID_For_PANDARoot)

- https://pandaatfair.githost.io/WEsmail/PandaRoot

# Available Multivariate packages:

- TMVA (ROOT).

- MLlib Apache Spark.

- Sci-Kit Learn.

- TensorFlow (Deep Learning).

- Keras (Deep Learning).

- PyTorch (Deep Learning).

- DL4J (Deep Learning, Java).

- R Implementations.

- …

JÜLICH
Forschungszentrum

# Conclusion and outlook:

- *Boosted Decision Trees (BDT)* showed good performance in classifying charged particles and outperformed the classical PID methods over the specified momentum range.

- *Deep Networks* also showed also good performance in the classifying task (needs more deep investigations).

- Interface between Python and PandaRoot by ZeroMQ.

- Write a release note of what have been done (*in progress …*).

**JÜLICH**
Forschungszentrum

# THANK YOU

**JÜLICH**
Forschungszentrum

# BACK UP

# Boosting:

1.  Draw a random subset of training samples $d_1$ without replacement from the training set $D$ to train a weak learner $C_1$.

2.  Draw second random training subset $d_2$ without replacement from the training set and add 50 percent of the samples that were previously misclassified to train a weak learner $C_2$.

3.  Find the training samples $d_3$ in the training set $D$ on which $C_1$ and $C_2$ disagree to train a third weak learner $C_3$.

4.  Combine the weak learners $C_1$, $C_2$, and $C_3$ via majority voting.

**JÜLICH**
Forschungszentrum