首先需要登录交大服务器

输入

ssh -Y inpactest@bl-1-1.physics.sjtu.edu.cn
password: inpac123456

# Step1 Training

1. 首先进入/Erec_tuto/tmva 路径下的 TMVARegression.C
   cd /Erec_tuto/tmva
   vi TMVARegression.C

2. 设定 MVA 中所使用的方法。1 表示打开，0 表示关闭。本次所使用的方法为 MLP 和
   BDTG 所以请确保 MLP 和 BDTG 方法打开。将程序调整为下图所示。

```
73   // --- Mutidimensional likelihood and Nearest-Neighbour methods
74   Use["PDERS"]        = 0;
75   Use["PDEFoam"]      = 0;
76   Use["KNN"]          = 0;
77   //
78   // --- Linear Discriminant Analysis
79   Use["LD"]          = 0;
80   //
81   // --- Function Discriminant analysis
82   Use["FDA_GA"]       = 0;
83   Use["FDA_MC"]       = 0;
84   Use["FDA_MT"]       = 0;
85   Use["FDA_GAMT"]     = 0;
86   //
87   // --- Neural Network
88   Use["MLP"]          = 1;
89   //
90   // --- Support Vector Machine
91   Use["SVM"]          = 0;
92   //
93   // --- Boosted Decision Trees
94   Use["BDT"]          = 0;
95   Use["BDTG"]         = 1;
96   // ----------------------------------------------------------
```

3. 创建 factory object. 之后你可以对不同的 factory 选择不同的 TMVA 的方法。如果
   程序用作回归，请将第一个选项调整为"TMVARegression"，用作分类请调整为
   "TMVAClassification"。本次重建能量用作回归，所以将程序调整至下图所示。

```
137    TMVA::Factory *factory = new TMVA::Factory( "TMVARegression", outputFile,
138                        "!V:!Silent:Color:DrawProgressBar" );
```

4. 定义输入变量，这些输入变量将被应用到 MVA 的 training 中。你也可以使用一些变
   量的组合类似于"3*var1/var2*abs(var3)"。下图为我们所使用到的 12 个输入变量。

请确保你的程序和下图相同。

-第一个选项是该变量在 root 文件中的名字。

-第二个选项是你希望在输出结果中显示的名字。

-第三个选项是该变量的单位。

-第四个选项是该变量的数据类型。"I"表示"Int"，"F"表示"Float"……

```
149    factory->AddVariable( "nHit", "nHit", "units", 'F' );
150    factory->AddVariable( "nHit1", "nHit1", "units", 'F' );
151    factory->AddVariable( "nHit2", "nHit2", "units", 'F' );
152    factory->AddVariable( "nHit3", "nHit3", "units", 'F' );
153    factory->AddVariable( "nHough", "nHough", "", 'I' );
154    factory->AddVariable( "nCluster", "nCluster", "", 'I' );
155    factory->AddVariable( "nTrack", "nTrack", "", 'I' );
156    factory->AddVariable( "nLayer", "nLayer", "", 'I' );
157    factory->AddVariable( "density", "Density", "", 'F' );
158    factory->AddVariable( "meanRadius", "meanRadius", "", 'F' );
159    factory->AddVariable( "nInteractingLayer", "Interlayer", "", 'I' );
160    factory->AddVariable( "begin", "begin", "", 'I' );
```

5. 定义作为 target 的变量。本次是为了重建能量。因此将 energy 作为程序的 target。

```
175    factory->AddTarget( "energy" );
```

6. 添加用作 TMVA training 的 MC 数据文件。在路径/Erec_tuto/trainingfile/root_training 下面，我们为你提供了四个 pi-的 root 文件，它们的区别仅在于数据量的不同，分别为 1k，5k，1w，5w。在本程序中我们添加的为数据量为 1w 的文件，但你也可以改变为其他文件，观察数据量不同所导致的结果的差异，当然数据量越大，程序运行的时间就越长。

```
183    TFile *input(0);
184    TString fname = "../trainingfile/root_training/pi-_1w.root";
```

7. 添加所使用的 root 文件中的 tree。

```
198    TTree *regTree = (TTree*)input->Get("tree");
```

8. 设置 TMVA 中各个方法的参数。本次所示用的方法是 MLP 和 BDTG，因此只需要设置这两个方法的参数即可。下图为本程序中 MLP 方法的设置选项。

```
268  if (Use["MLP"])
269    factory->BookMethod( TMVA::Types::kMLP, "MLP", "!H:!V:VarTransform=Norm:NeuronType=tanh:NCycles=
5000:HiddenLayers=10,2:TestRate=6:TrainingMethod=BFGS:Sampling=0.3:SamplingEpoch=0.8:
ConvergenceImprove=1e-6:ConvergenceTests=15:!UseRegulator:LearningRate=0.001" );
```

下表中为 MLP 方法中各个参数的意义和解释，你可以根据下表中的信息，对上图中的各个参数进行一定的修改和调试。你也可以参照 TMVA_Users guide 进行进一步了解。

| Option | Array | Default | Predefined Values | Description |
|---|---|---|---|---|
| NCycles | − | 500 | − | Number of training cycles |
| HiddenLayers | − | N,N-1 | − | Specification of hidden layer architecture |
| NeuronType | − | sigmoid | linear, sigmoid, tanh, radial | Neuron activation function type |
| NeuronInputType | − | sum | sum, sqsum, abssum | Neuron input function type |
| TrainingMethod | − | BP | BP, GA, BFGS | Train with Back-Propagation (BP), BFGS Algorithm (BFGS), or Genetic Algorithm (GA - slower and worse) |
| LearningRate | − | 0.02 | − | ANN learning rate parameter |
| DecayRate | − | 0.01 | − | Decay rate for learning parameter |
| TestRate | − | 10 | − | Test for overtraining performed at each #th epochs |
| Sampling | − | 1 | − | Only 'Sampling' (randomly selected) events are trained each epoch |
| SamplingEpoch | − | 1 | − | Sampling is used for the first 'SamplingEpoch' epochs, afterwards, all events are taken for training |
| SamplingImportance | − | 1 | − | The sampling weights of events in epochs which successful (worse estimator than before) are multiplied with SamplingImportance, else they are divided. |

| Option | Array | Default | Predefined Values | Description |
|---|---|---|---|---|
| SamplingTraining | − | True | − | The training sample is sampled |
| SamplingTesting | − | False | − | The testing sample is sampled |
| ResetStep | − | 50 | − | How often BFGS should reset history |
| Tau | − | 3 | − | LineSearch size step |
| BPMode | − | sequential | sequential, batch | Back-propagation learning mode: sequential or batch |
| BatchSize | − | -1 | − | Batch size: number of events/batch, only set if in Batch Mode, -1 for BatchSize=number_of_events |
| ConvergenceImprove | − | 0 | − | Minimum improvement which counts as improvement (<0 means automatic convergence check is turned off) |
| ConvergenceTests | − | -1 | − | Number of steps (without improvement) required for convergence (<0 means automatic convergence check is turned off) |

下图为本程序中 BDTG 方法的设置选项。

```
280    if (Use["BDTG"])
281      factory->BookMethod( TMVA::Types::kBDT, "BDTG",
282             "!H:!V:NTrees=2000::BoostType=Grad:Shrinkage=0.1:UseBaggedBoost:BaggedSampleFraction=0.5:
    nCuts=20:MaxDepth=3:MaxDepth=4" );
```
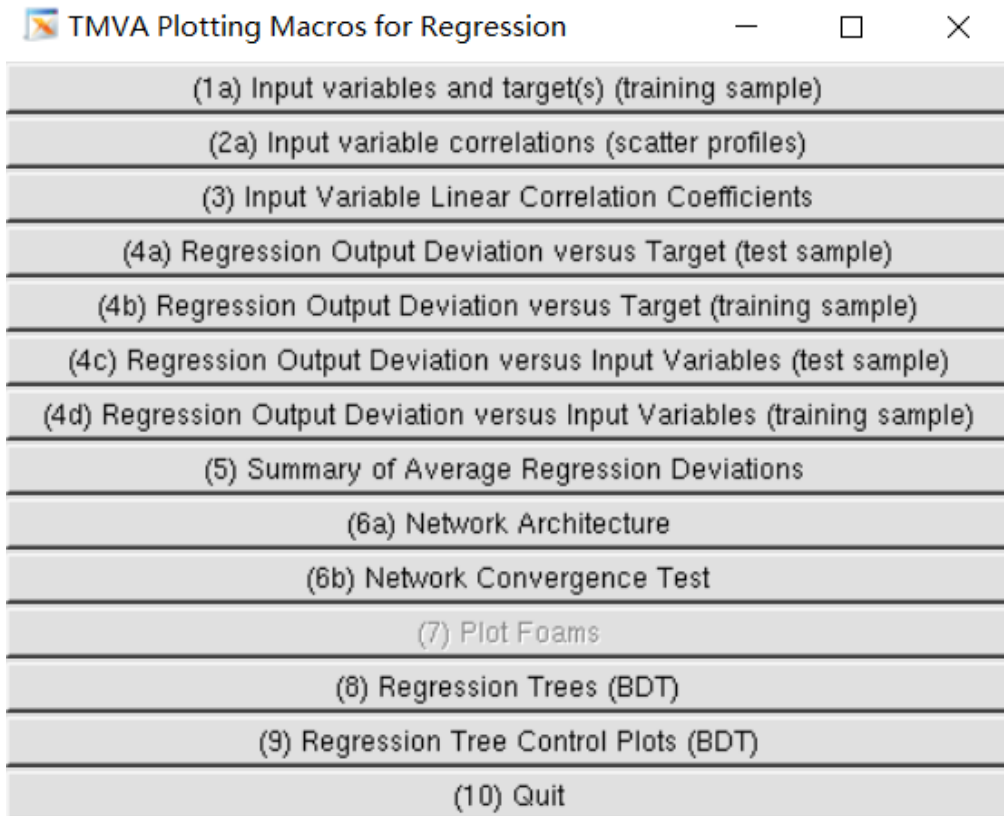
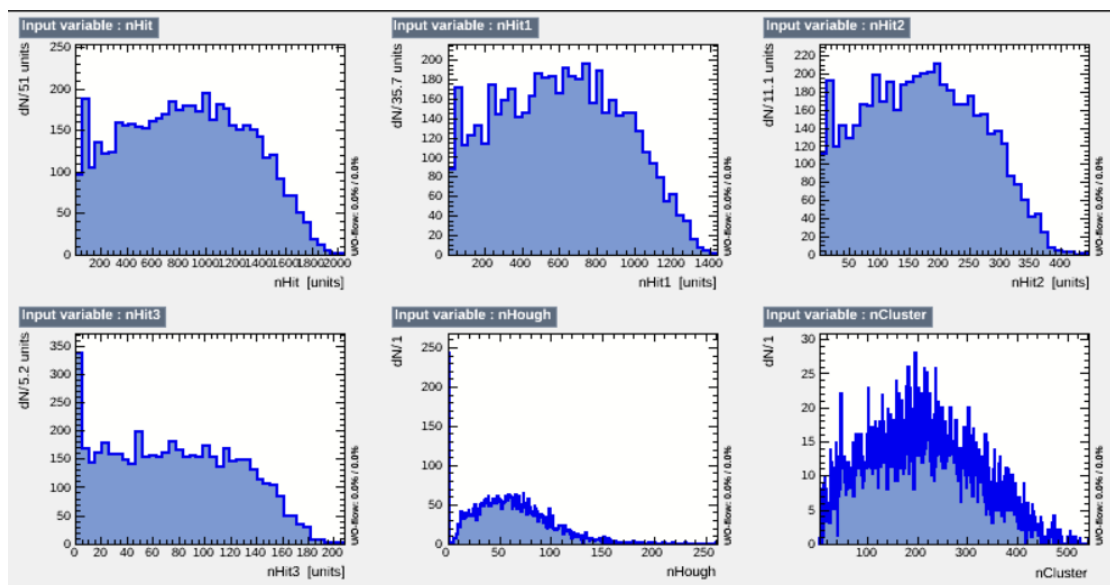下表中为 BDTG 方法中各个参数的意义和解释，你可以根据下表中的信息，对上图中的各个参数进行一定的修改和调试。你也可以参照 TMVA_Users guide 进行进一步了解。

| Option | Array | Default | Predefined Values | Description |
| --- | --- | --- | --- | --- |
| SeparationType | — | GiniIndex | CrossEntropy, GiniIndex, GiniIndexWithLaplace, MisClassificationError, SDivSqrtSPlusB, RegressionVariance | Separation criterion for node splitting |
| nEventsMin | — | max(20,NEvtsTrain/NVar$^2$/10) | | Minimum number of events required in a leaf node (default uses given formula) |
| nCuts | — | 20 | — | Number of steps during node cut optimisation |
| PruneStrength | — | -1 | — | Pruning strength |
| PruneMethod | — | CostComplexity | NoPruning, ExpectedError, CostComplexity | Method used for pruning (removal) of statistically insignificant branches |
| PruneBeforeBoost | — | False | — | Flag to prune the tree before applying boosting algorithm |
| PruningValFraction | — | 0.5 | — | Fraction of events to use for optimizing automatic pruning. |
| NNodesMax | — | 100000 | — | Max number of nodes in tree |
| MaxDepth | — | 100000 | — | Max depth of the decision tree allowed |

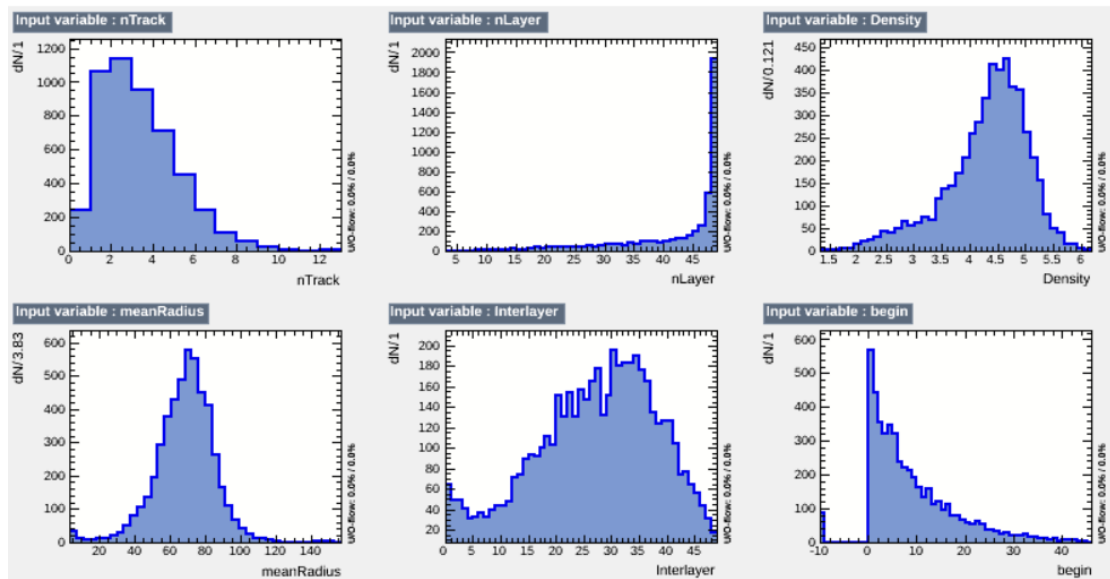| Option | Array | Default | Predefined Values | Description |
| --- | --- | --- | --- | --- |
| NTrees | — | 200 | — | Number of trees in the forest |
| BoostType | — | AdaBoost | AdaBoost, Bagging, RegBoost, AdaBoostR2, Grad | Boosting type for the trees in the forest |
| AdaBoostR2Loss | — | Quadratic | Linear, Quadratic, Exponential | Loss type used in AdaBoostR2 |
| UseBaggedGrad | — | False | — | Use only a random subsample of all events for growing the trees in each iteration. (Only valid for GradBoost) |
| GradBaggingFraction | — | 0.6 | — | Defines the fraction of events to be used in each iteration when UseBaggedGrad=kTRUE. |
| Shrinkage | — | 1 | — | Learning rate for GradBoost algorithm |
| AdaBoostBeta | — | 1 | — | Parameter for AdaBoost algorithm |
| UseRandomisedTrees | — | False | — | Choose at each node splitting a random set of variables |
| UseNvars | — | 4 | — | Number of variables used if randomised tree option is chosen |
| UseNTrainEvent | — | N | — | Number of Training events used in each tree building if randomised tree option is chosen |
| UseWeightedTrees | — | True | — | Use weighted trees or simple average in classification from the forest |
| UseYesNoLeaf | — | True | — | Use Sig or Bkg categories, or the purity=S/(S+B) as classification of the leaf node |
| NodePurityLimit | — | 0.5 | — | In boosting/pruning, nodes with purity > NodePurityLimit are signal; background otherwise. |

9. 到此为止能量重建 training 部分的程序已设定完毕，退出程序，运行程序即可。
   root TMVARegression.C

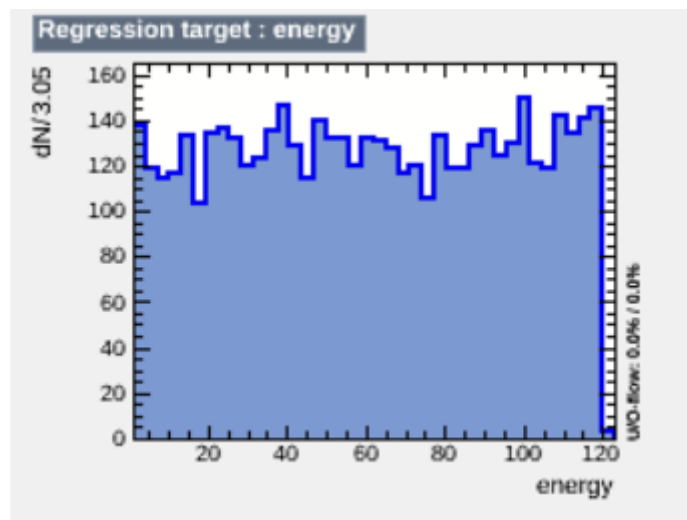10. 输出结果。当程序运行结束之后会生成一个控制面板。你可以点击相应的选

    项来查看不同的输出结果。你需要主要了解的选项如下。
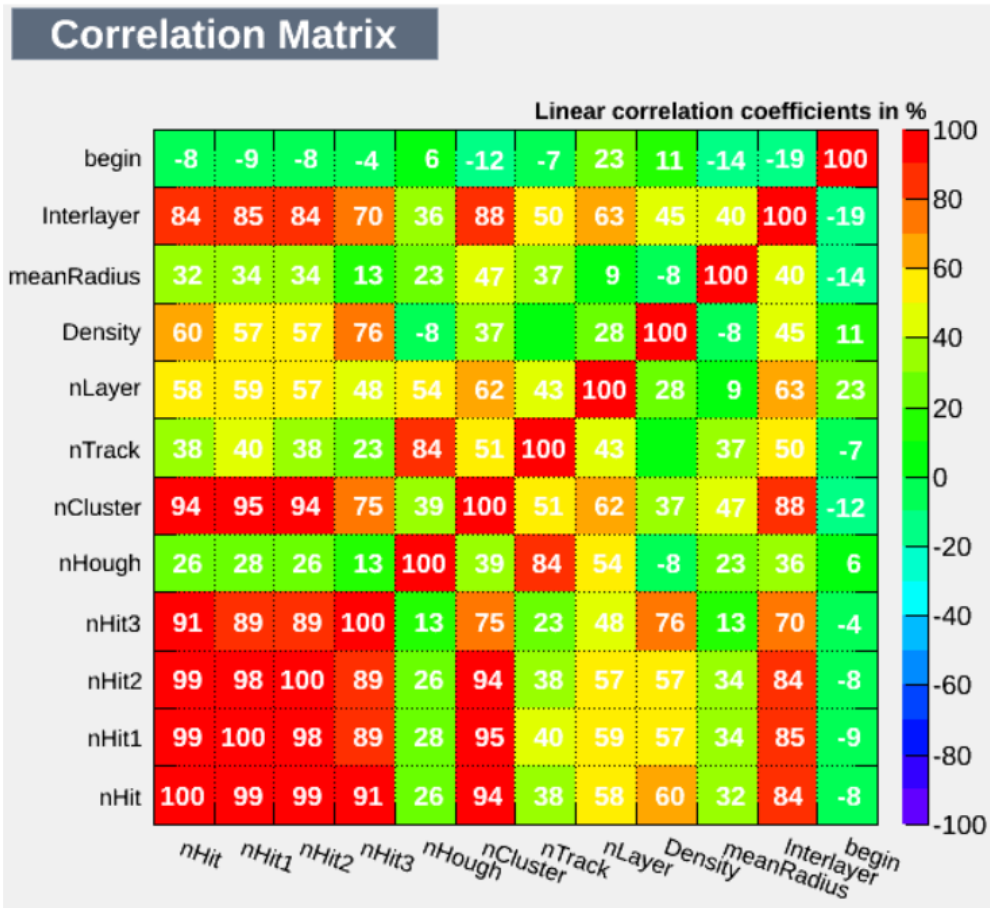


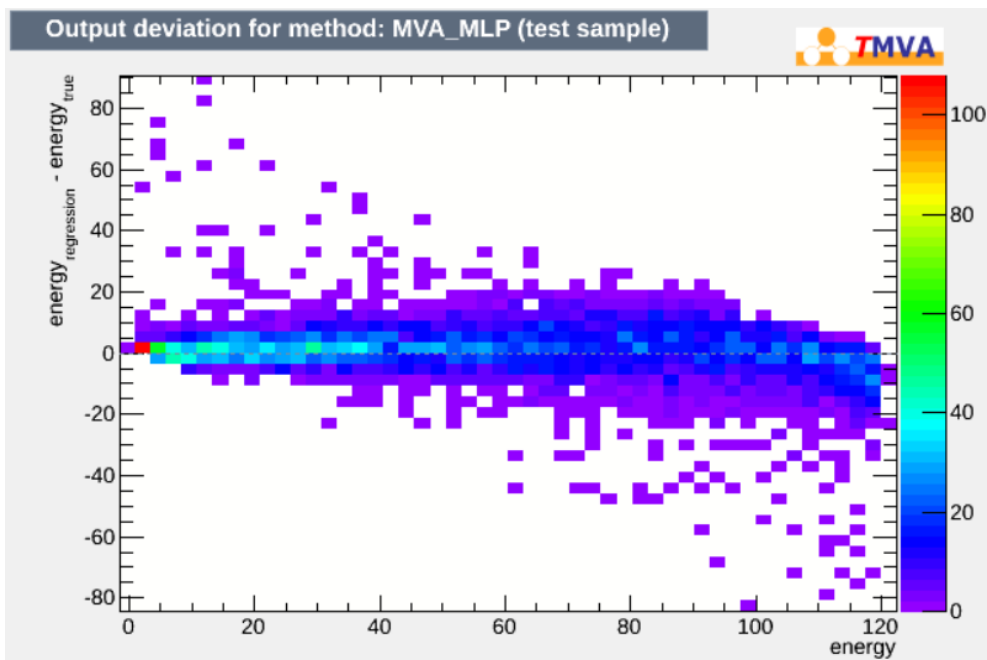1) 1a 表示输入 variables 和 target 的分布。
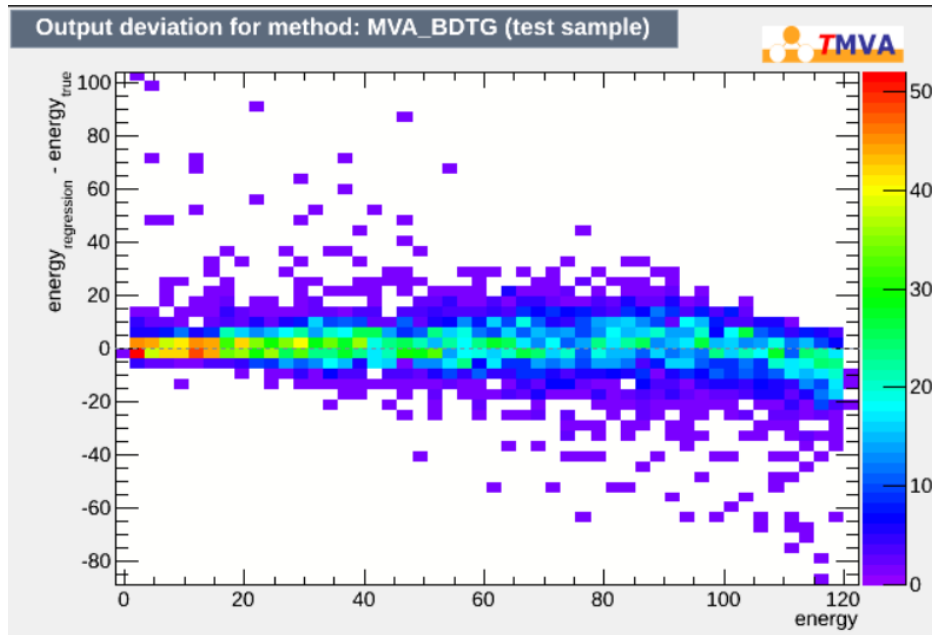


（input variables）

（input variables）



（target）

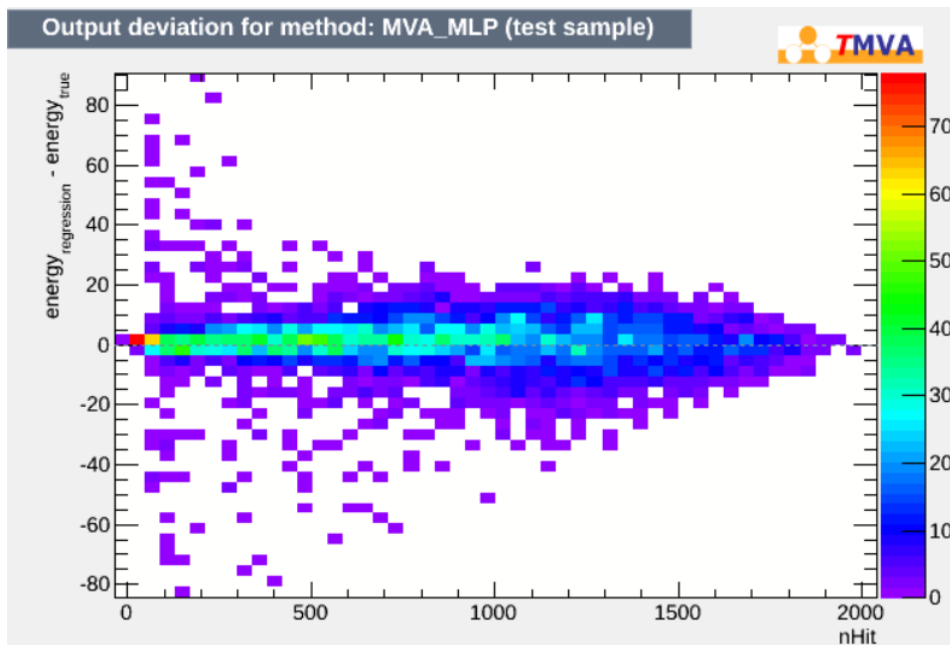2) 3a 为输入各变量之间的相关度

3) 4a 为能量重建 test 的结果。纵轴为重建能量-真实能量，横轴为真实能量。



（MLP 方法 test 结果）

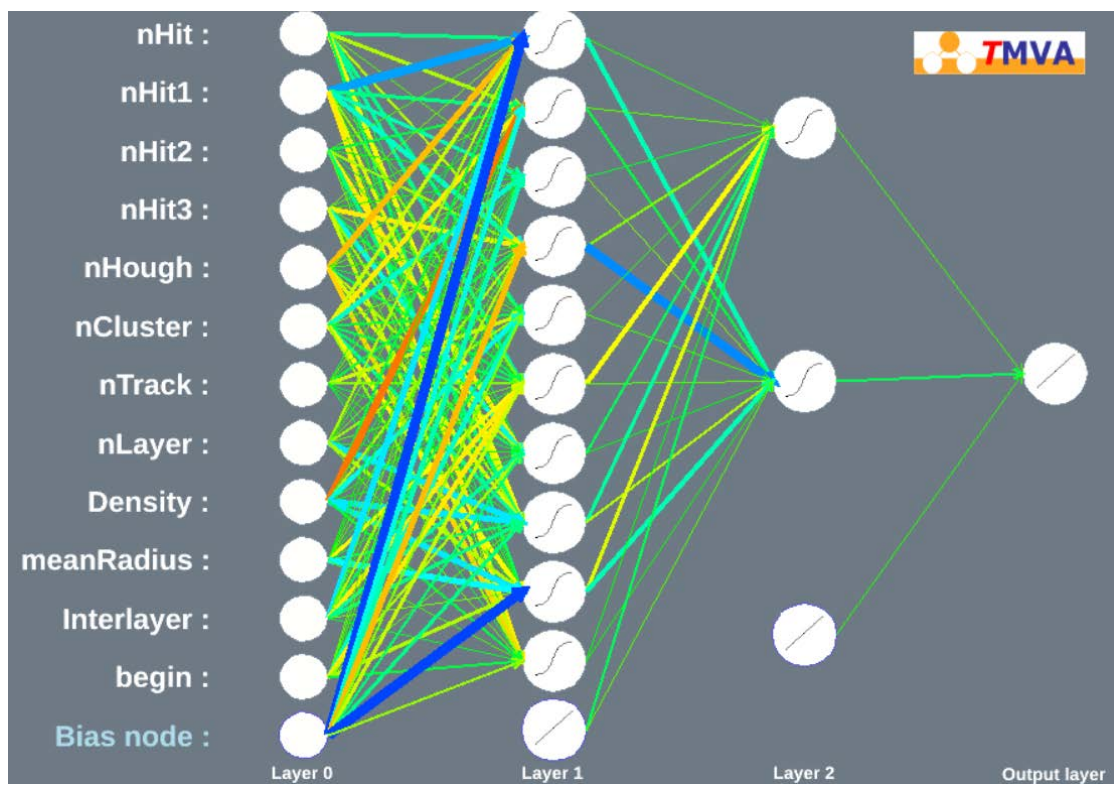（BDTG 方法 test 结果）

4) 4c 也是能量重建 test 的结果，但是横轴更换为各个不同的变量，你可以通过 4c 展示 fig 中的分布判断不同变量对能量重建的影响。例如：nHit 对能量重建的影响，如图所示：



5) 6a 为 MLP 的结构图，你可以在步骤 8 中调整这个结构。

# Step2 Application

1. 首先进入/Erec_tuto/tmva 路径下的 TMVARegressionApplication.C
   cd /Erec_tuto/tmva
   vi TMVARegressionApplication.C
2. 设定 MVA 中所使用的方法。1 表示打开，0 表示关闭。本次所使用的方法为 MLP 和 BDTG 所以请确保 MLP 和 BDTG 方法打开。将程序调整为下图所示。

```
41          // --- Mutidimensional likelihood and Nearest-Neighbour methods
42          Use["PDERS"]            = 0;
43          Use["PDEFoam"]          = 0;
44          Use["KNN"]              = 0;
45          //
46          // --- Linear Discriminant Analysis
47          Use["LD"]               = 0;
48          //
49          // --- Function Discriminant analysis
50          Use["FDA_GA"]           = 0;
51          Use["FDA_MC"]           = 0;
52          Use["FDA_MT"]           = 0;
53          Use["FDA_GAMT"]         = 0;
54          //
55          // --- Neural Network
56          Use["MLP"]              = 1;
57          //
58          // --- Support Vector Machine
59          Use["SVM"]              = 0;
60          //
61          // --- Boosted Decision Trees
62          Use["BDT"]              = 0;
63          Use["BDTG"]             = 1;
64          //
```

3. 添加变量，这些变量<mark>必须</mark>和 TMVARegression.C 中的变量相同

```
95          Float_t nHit1, nHit2, nHit3 ,nHit ,nHough ,nClusteer ,nTrack ,nLayer  ,nInter ;
96          Float_t nbegin,ndensity ,nmeanRadius;
97          reader->AddVariable( "nHit", &nHit );
98          reader->AddVariable( "nHit1", &nHit1 );
99          reader->AddVariable( "nHit2", &nHit2 );
100         reader->AddVariable( "nHit3", &nHit3 );
101         reader->AddVariable( "nHough", &nHough);
102         reader->AddVariable( "nCluster", &nClusteer);
103         reader->AddVariable( "nTrack", &nTrack);
104         reader->AddVariable( "nLayer", &nLayer);
105         reader->AddVariable( "density", &ndensity);
106         reader->AddVariable( "meanRadius", &nmeanRadius);
107         reader->AddVariable( "nInteractingLayer", &nInter);
108         reader->AddVariable( "begin", &nbegin);
```

4. 设置重建结束之后输出重建能量分布的参数。其中 100,0,20 分别表示 bin 的数量，histgram 的能量范围，我们将其设定为 0,20。当你重建的能量在此区域之外时可以适当调整该值。

```
132         for (std::map<std::string,int>::iterator it = Use.begin(); it != Use.end(); it++) {
133             TH1* h = new TH1F( it->first.c_str(), TString(it->first) + " method",100 ,0, 20);
134             if (it->second) hists[++nhists] = h;
```

5. 添加需要重建能量的 MC/data samples 的路径，我们在路径 Erec_tuto/trainingfile/root_app 下为你准备了 8 个 pi-的 MC root 文件，这些 root 文件的区别在于它们拥有不同的分立的能量。分别是 10,20,30……80GeV，你可以修改下面这一行代码，尝试对不同能量的 root 文件进行能量重建。

```
143         TString fname = TString::Format("../trainingfile/root_app/pi-_10GeV.root");
```

6. 读取 root 中的 tree 和 variables，每一个 variable 都需要与之前所填写的一一对应。

```
163        TTree* theTree = (TTree*)input->Get("tree");
164        Int_t Nhit1,Nhit2,Nhit3,Nhit,Nhough,NClusteers,ntrack,Nlayer,NInter;
165        Double_t nDensity,nRadius,nBegin;
166        std::cout << "--- Select signal sample" << std::endl;
167        theTree->SetBranchAddress( "nHit", &Nhit );
168        theTree->SetBranchAddress( "nHit1", &Nhit1 );
169        theTree->SetBranchAddress( "nHit2", &Nhit2 );
170        theTree->SetBranchAddress( "nHit3", &Nhit3 );
171        theTree->SetBranchAddress( "nHough", &Nhough);
172        theTree->SetBranchAddress( "nCluster", &NClusteers);
173        theTree->SetBranchAddress( "nTrack", &ntrack);
174        theTree->SetBranchAddress( "nLayer", &Nlayer);
175        theTree->SetBranchAddress( "density", &nDensity);
176        theTree->SetBranchAddress( "meanRadius", &nRadius);
177        theTree->SetBranchAddress( "nInteractingLayer", &NInter);
178        theTree->SetBranchAddress( "begin", &nBegin);
```

7. 定义一个新的 root 文件，并将重建能量的分布存储到这个 root 文件中。
```
238                TString foutName = TString::Format("TMVARegApp_10GeV.root");
```

8. 退出文件，并运行该文件即可用之前 training 过的软件包对新的数据进行能量重建。
   root TMVARegressionApplication.C

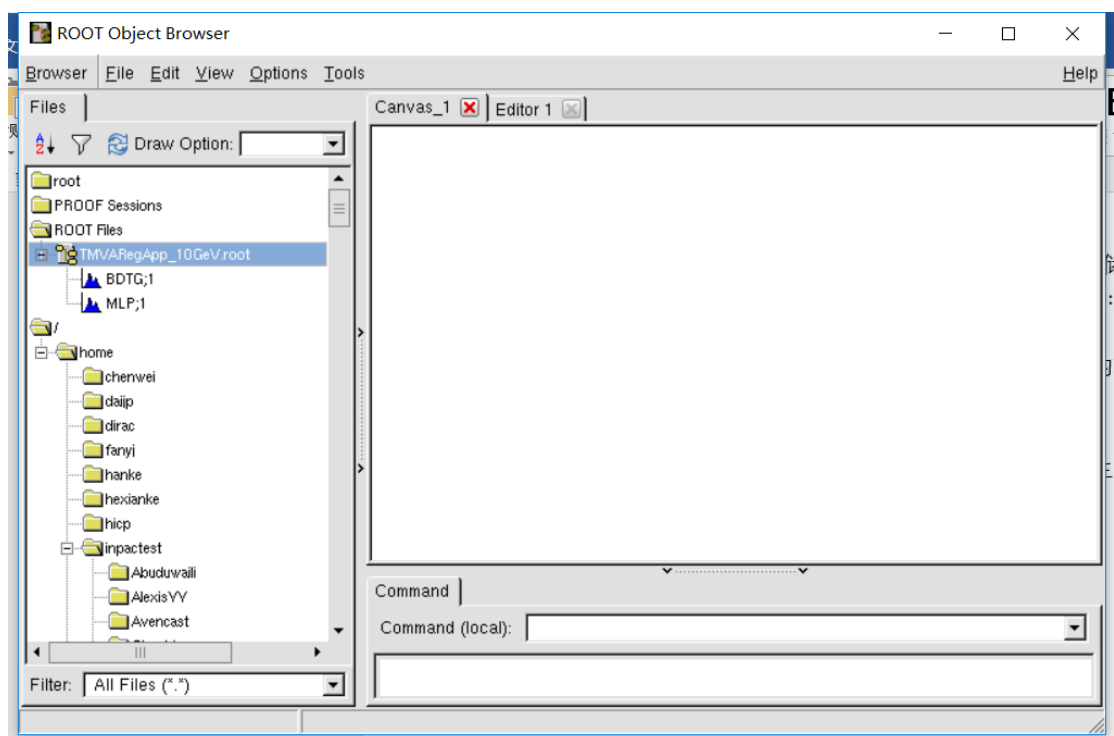9. 运行完第 7 步之后会产生一个新的 tmva 目录下产生一个新的 root 文件，也就是第 6 步
   中所定义的"TMVARegApp_10GeV.root"。
   打开这个 root 文件。
   root TMVARegressionApp_10GeV.C
   输入
   root [1] TBrowser a
   会进入浏览界面

双击"BDTG；1"和"MLP；1"两个图标即可得到能量重建的结果。



BDTG method

| BDTG | |
|---|---|
| Entries | 20000 |
| Mean | 10.32 |
| RMS | 2.606 |

MLP method

| MLP | |
|---|---|
| Entries | 20000 |
| Mean | 10.29 |
| RMS | 3.001 |