# Ejets Update 28/05/2020

# BDT (TMVA) Implementation for Selection

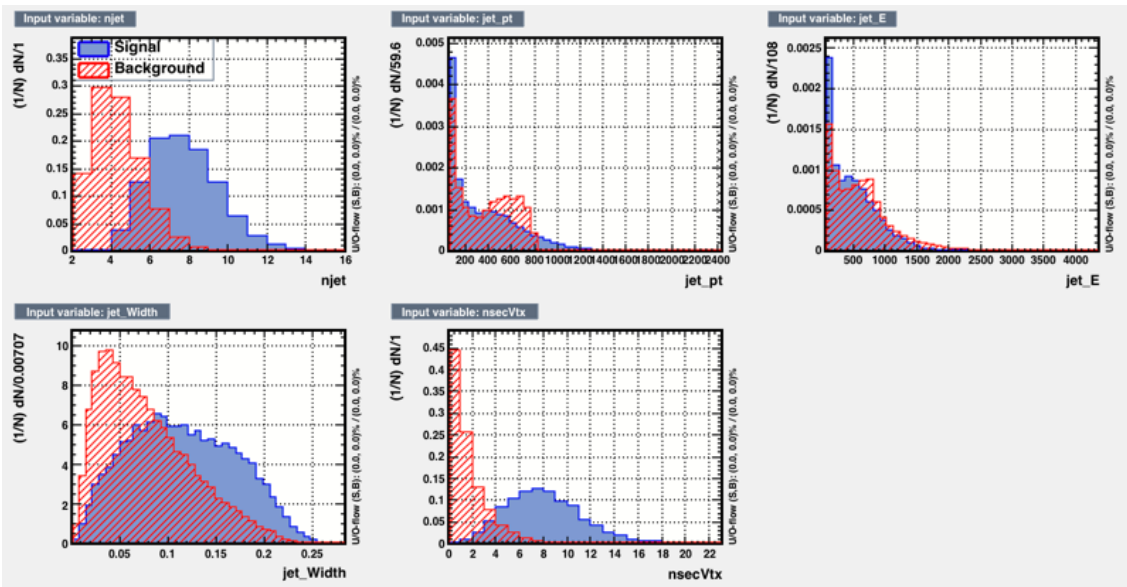| Option | Array | Default | Predefined Values | Description |
|---|---|---|---|---|
| NTrees | — | 800 | — | Number of trees in the forest |
| MaxDepth | — | 3 | — | Max depth of the decision tree allowed |
| MinNodeSize | — | 5% | — | Minimum percentage of training events required in a leaf node (default: Classification: 5%, Regression: 0.2%) |
| nCuts | — | 20 | — | Number of grid points in variable range used in finding optimal cut in node splitting |
| BoostType | — | AdaBoost | AdaBoost, RealAdaBoost, Bagging, AdaBoostR2, Grad | Boosting type for the trees in the forest (note: AdaCost is still experimental) |
| AdaBoostR2Loss | — | Quadratic | Linear, Quadratic, Exponential | Type of Loss function in AdaBoostR2 |
| UseBaggedGrad | — | False | — | Use only a random subsample of all events for growing the trees in each iteration. (Only valid for GradBoost) |
| Shrinkage | — | 1 | — | Learning rate for GradBoost algorithm |
| AdaBoostBeta | — | 0.5 | — | Learning rate for AdaBoost algorithm |
| UseRandomisedTrees | — | False | — | Determine at each node splitting the cut variable only as the best out of a random subset of variables (like in RandomForests) |
| UseNvars | — | 2 | — | Size of the subset of variables used with RandomisedTree option |
| UsePoissonNvars | — | True | — | Interpret UseNvars not as fixed number but as mean of a Possion distribution in each split with RandomisedTree option |
| BaggedSampleFraction | — | 0.6 | — | Relative size of bagged event sample to original size of the data sample (used whenever bagging is used (i.e. UseBaggedGrad, Bagging,) |
| UseYesNoLeaf | — | True | — | Use Sig or Bkg categories, or the purity=S/(S+B) as classification of the leaf node -> Real-AdaBoost |

# Summary

- have implemented the ROOT TMVA BDT tool on signal/background MC

- used all events from files with no cuts

- did not apply event weighting

- treated all variables as 'event level' variables
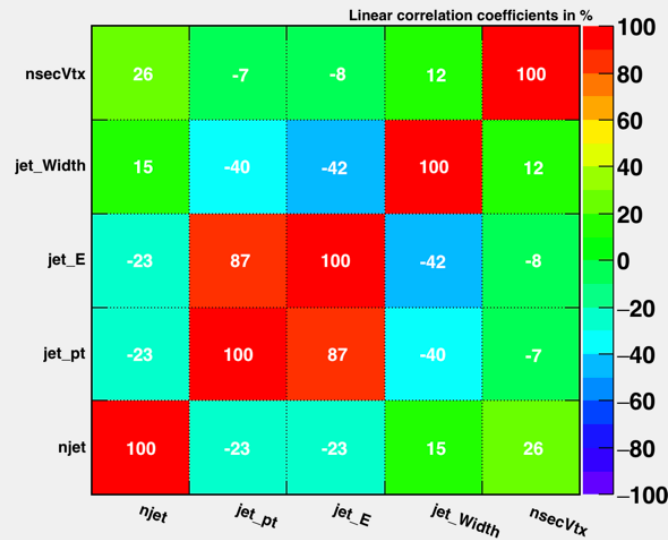
- used default BDT settings

- ran over cases:

| | | |
|---|---|---|
| ModelA_1400_20 | ModelB_1400_20 | ModelE_1400_75 |
| ModelA_1000_150 | ModelB_1000_5 | ModelE_1000_150 |
| ModelA_600_1 | ModelB_600_300 | ModelE_600_0p5 |

- Used variables:  **njet   jet_pt   jet_E   jet_Width   nsecVtx**

**This is not a real selection/analysis...just proof-of-principle test...**

# ModelA_1400_20

# ModelA_600_1

# ModelB_600_300

# ModelE 1000 150

# Now try removing variables...again with Model A 1400 20



Drop:     1. njet
          2. jet_pt
          3. jet_E
          4. jet_Width
          5. nsecVtx

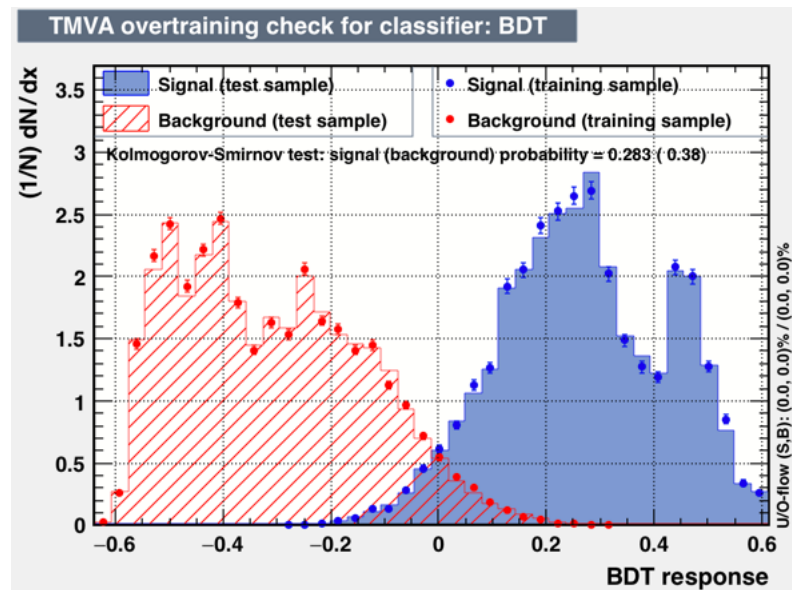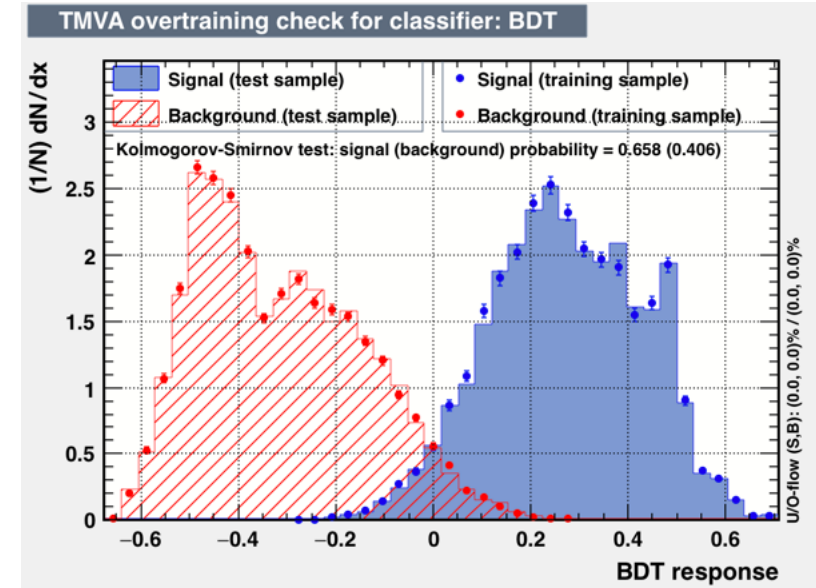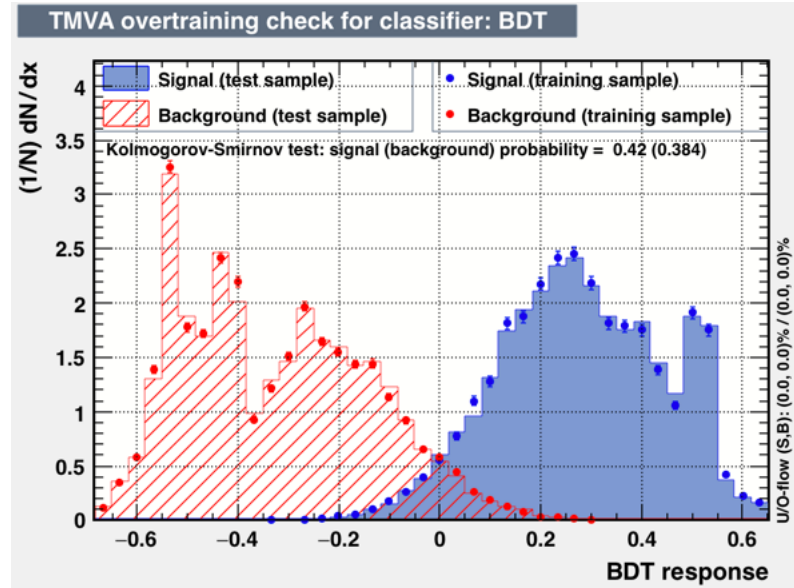# Now try removing variables...again with Model A 1400 20



Drop:
1. njet
2. jet_pt
3. jet_E
4. jet_Width
5. nsecVtx

# Now try removing variables...again with Model A 1400 20
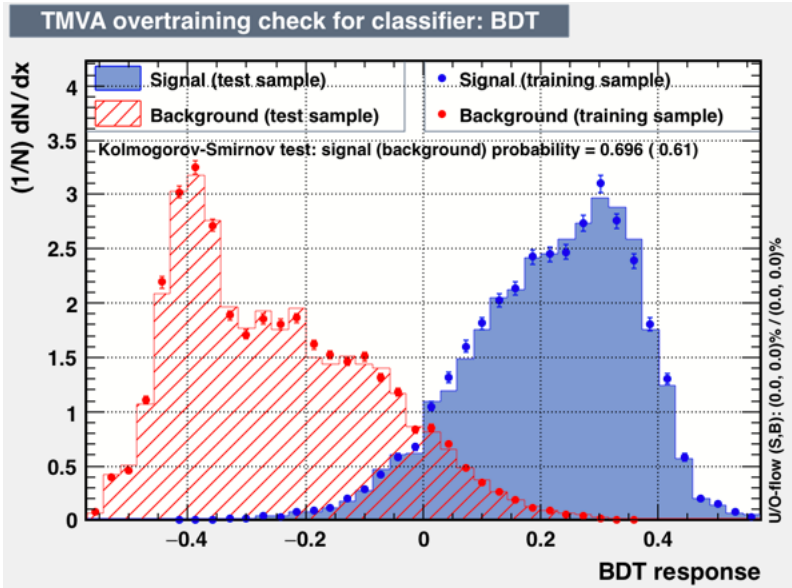


Drop:
1. njet
2. jet_pt
3. jet_E
4. jet_Width
5. nsecVtx

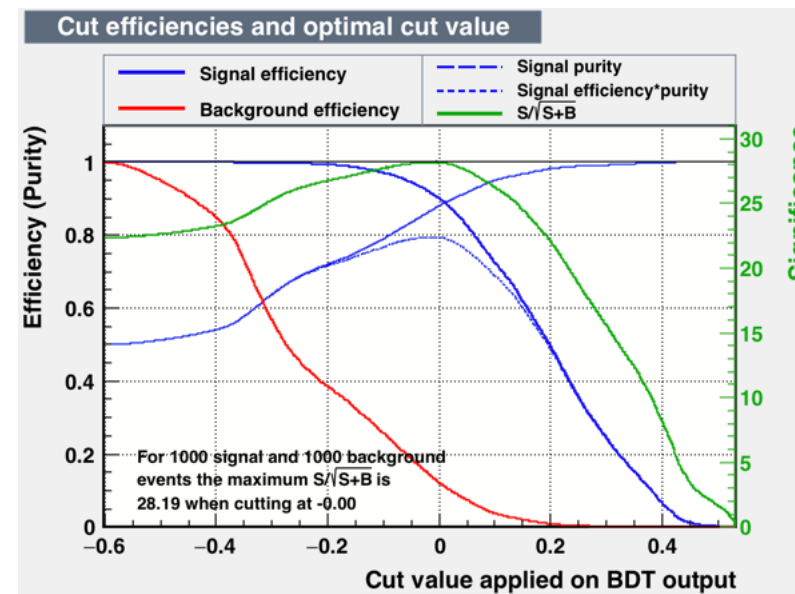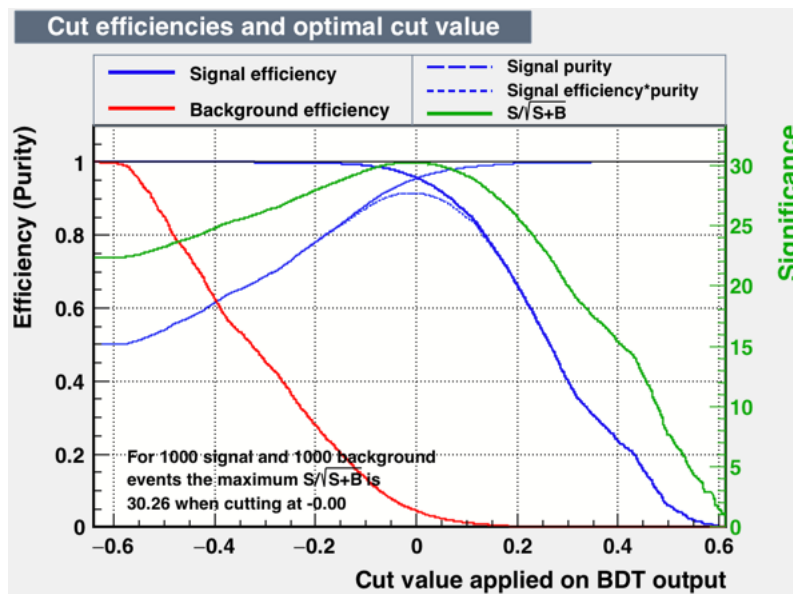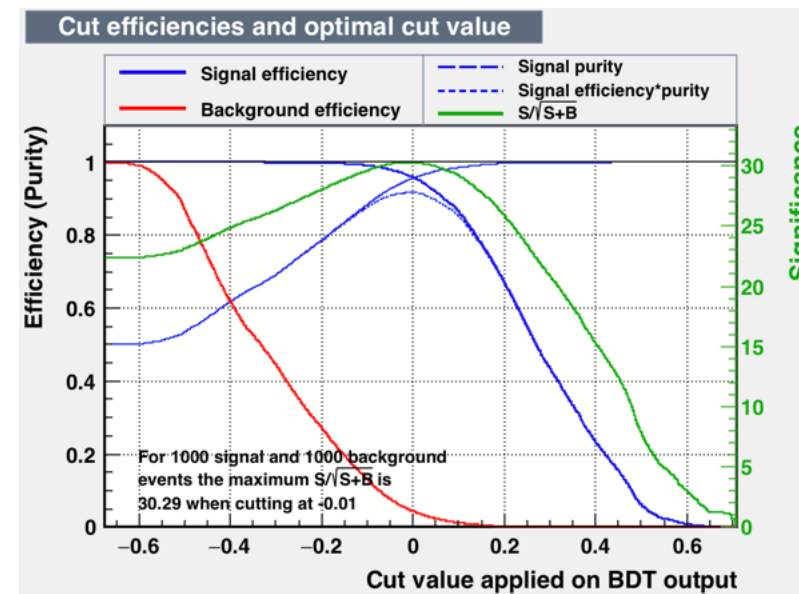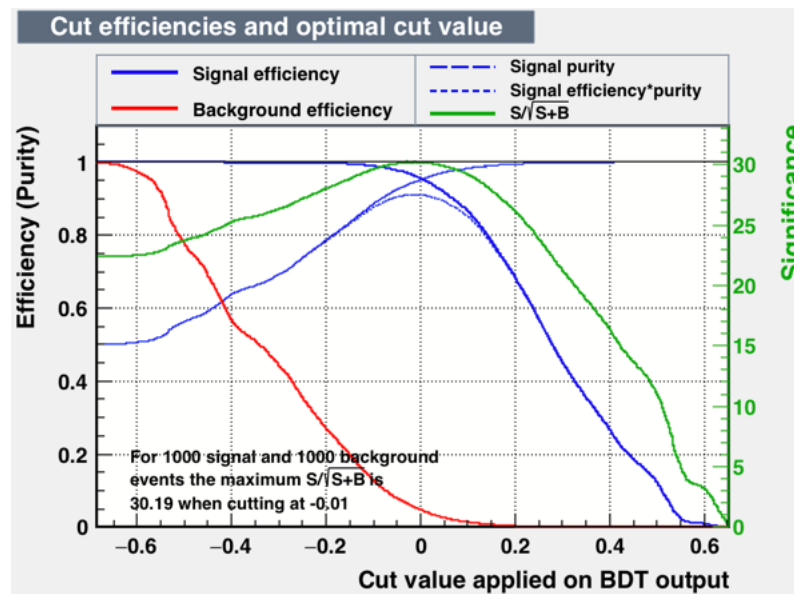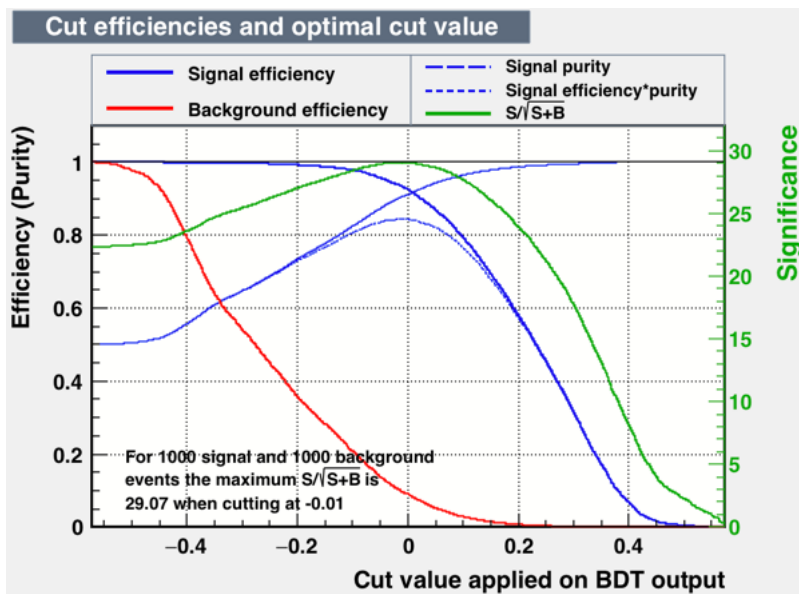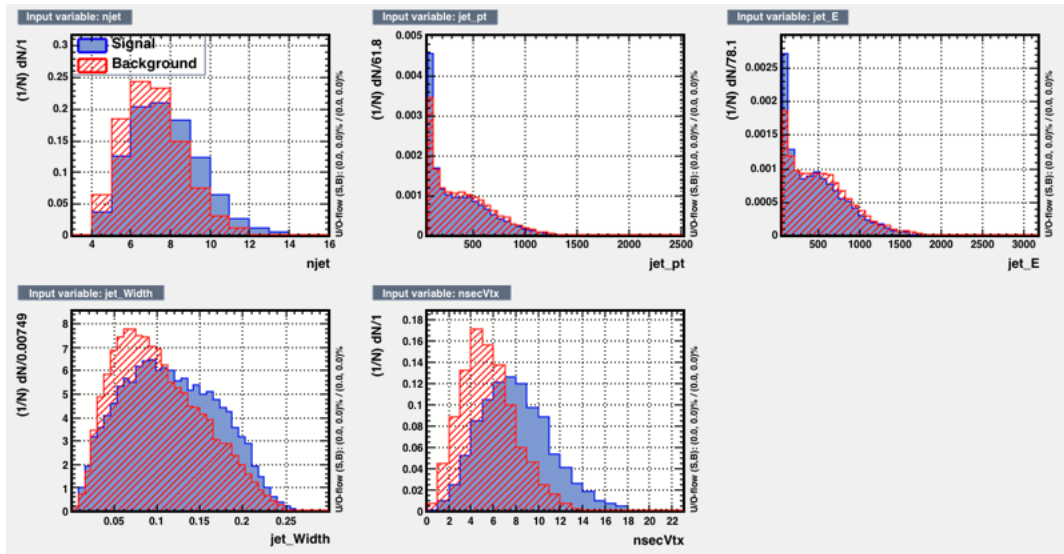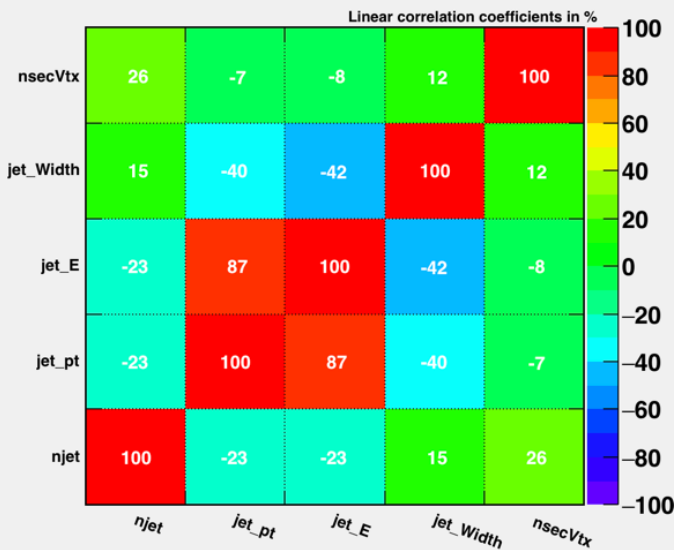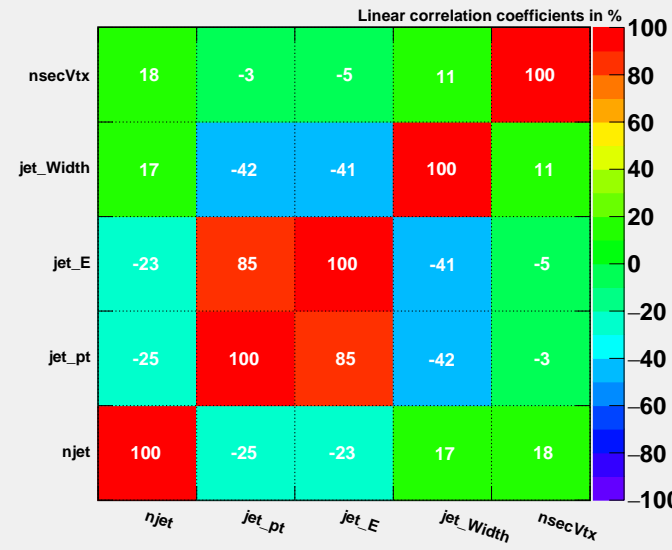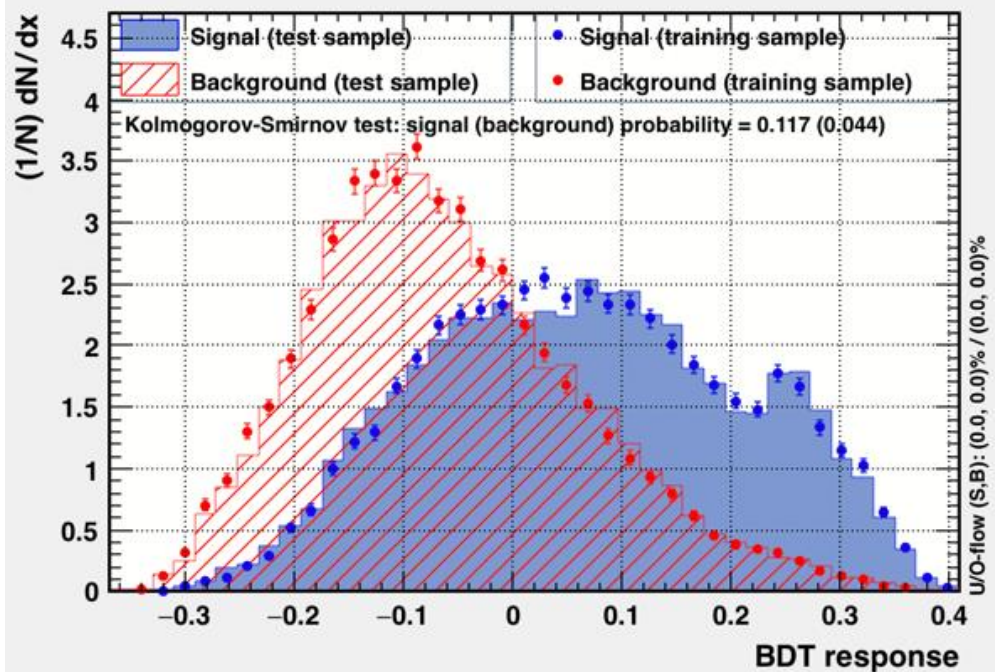# Try A_1400_20 against B_1400_20

# Comments

- wrong to treat variables as event variables…just a test
    - plan to have sub-jet-level variables, jet-variables, event variables…multi-step approach
    - also need to apply pre-selection cuts and correct event weightings
    - MC statistics are on the low side

- use as a tool to test sensitivity of variables…but need to be aware of modelling fidelity
    - need to connect to sensitivity as benchmark…maybe S/sqrt(S+B) OK though?

- can check variable correlations

- can isolate events with 'distinct' BDT space features…should map to reconstructed space features

- easy to try alternate TMVA implementations and compare (e.g. NN, Fischer…)

- how to test systematics?  re-training with shifted distributions…but not necessarily easy to do..

- model dependence of selection is a challenge

- **I think this can be useful…but needs the implementation needs to go up in sophistication a few orders of magnitude…**