### **Machine learning**

with a high energy physicist bias

Yann Coadou

**CPPM** Marseille



### Online edition, 7 September 2021









#### IU" IDPASL SCHOOL

Nazaré, Portugal 06 - 16 September, 2021



### Introduction **Optimal discrimination Machine learning** Quadratic and linear discriminants (Boosted) Decision trees Neural networks **Deep neural networks** Machine learning and particle physics Conclusion References

Physics at future colliders lorgen D'Hondt

Higgs and Multi-Higgs TBD

Data acquisition in present and future experiments@CERN Sra Gaspar

Dark Matter in Coliders Benjamin Fuks

(Heavy) Flavour Physics in the precision ERA Stephane Montell

CERN, particle physics outreach Ana Godinho

The Origins of the Highest Energy Particles in Nature Alan Watson

Organizing Committee: M. Pimenta, N. Çastro , R. Gonçalo , P. Assis, I. Lopes , A. de Angelis Secretariat N. Antinos

https://indico.lip.pt/event/643/

FCT





### Typical problems in HEP

- Classification of objects
- Signal enhancement relative to background
- Regression: best estimation of a parameter
  - lepton energy,  $E_{\rm T}^{\rm miss}$  value, invariant mass, etc.

#### Discrimination of signal from background in HEP

- Event level (Higgs searches, ...)
- Cone level (tau-vs-jet reconstruction, ...)
- Lifetime and flavour tagging (*b*-tagging, ...)
- Track level (particle identification, ...)
- Cell level (energy deposit from hard scatter/pileup/noise, ...)



#### Input information from various sources

- Kinematic variables (masses, momenta, decay angles, ...)
- Event properties (jet multiplicity, sum of charges, brightness ...)
- Event shape (sphericity, aplanarity, ...)
- Detector response (silicon hits, dE/dx, Cherenkov angle, shower profiles, muon hits, ...)

### Most data are (highly) multidimensional

- Use dependencies between  $x = \{x_1, \cdots, x_n\}$  discriminating variables
- Approximate this *n*-dimensional space with a function f(x) capturing the essential features
- f is a multivariate discriminant
- For most of these lectures, use binary classification:
  - an object belongs to one class (e.g. signal) if f(x) > q, where q is some threshold,
  - and to another class (e.g. background) if  $f(x) \leq q$







• Where to place a cut x<sub>0</sub> on variable x?



• Optimal choice: minimum misclassification cost at decision boundary  $x = x_0$ 

# Poptimal discrimination: Bayes limit



### Bayes discriminant

$$BD = \frac{C_B}{C_S} = \frac{p(x|S)}{p(x|B)} \times \frac{p(S)}{p(B)} \qquad \begin{array}{c} C_B = \text{background contamination} \\ C_S = \text{signal loss} \end{array}$$
  
• From Bayes theorem  $(p(A|B)p(B) = p(B|A)p(A))$  and sum of probabilities  $(p(S|x) + p(B|x) = 1)$ :  
 $p(S|x) = \frac{BD}{1 + BD}$ 

#### **Bayes limit**

- p(S|x) = BD/(1 + BD) is what should be achieved to minimise cost of misclassification, reaching classification with fewest mistakes
- Fixing relative cost of background contamination and signal loss  $q = C_B/(C_S + C_B)$ , q = p(S|x) defines decision boundary:
  - signal-rich if  $p(S|x) \ge q$
  - background-rich if p(S|x) < q
- Any function that approximates conditional class probability p(S|x) with negligible error reaches the Bayes limit



### How to construct p(S|x)?

- k = p(S)/p(B) typically unknown
- Problem: p(S|x) depends on k!
- Solution: it's not a problem...
- Define a multivariate discriminant:

$$D(x) = \frac{s(x)}{s(x) + b(x)} = \frac{p(x|S)}{p(x|S) + p(x|B)}$$

Now:

$$p(S|x) = \frac{D(x)}{D(x) + (1 - D(x))/k}$$

• Cutting on D(x) is equivalent to cutting on p(S|x), implying a corresponding (unknown) cut on p(S|x)

# Reachine learning: learning from examples

### Several types of problems

- Classification/decision:
  - signal or background
  - type la supernova or not
  - will pay his/her credit back on time or not
- Regression: estimating a parameter value (energy of a particle, brightness of a supernova, ...) [mostly ignored in these lectures]
- Clustering (cluster analysis):
  - in exploratory data mining, finding features

#### Our goal

- Teach a machine to learn the discriminant f(x) using examples from a training dataset
- Be careful to not learn too much the properties of the training sample
  - no need to memorise the training sample
  - instead, interested in getting the right answer for new events  $\Rightarrow$  generalisation ability



# Rachine learning: training

### Supervised learning



- Training events are labelled: N examples (x, y)<sub>1</sub>, (x, y)<sub>2</sub>, ..., (x, y)<sub>N</sub> of (discriminating) feature variables x and class labels y
- The learner uses example classes to know how good it is doing

# Rachine learning: training

### Supervised learning



- Training events are labelled: N examples  $(x, y)_1, (x, y)_2, \dots, (x, y)_N$ of (discriminating) feature variables x and class labels y
- The learner uses example classes to know how good it is doing

#### **Reinforcement learning**

- Instead of labels, some sort of reward system (e.g. game score)
- Goal: maximise future payoff by optimising decision policy
- May not even "learn" anything from data, but remembers what triggers reward or punishment

# Rachine learning: training

### Supervised learning



- Training events are labelled: N examples (x, y)<sub>1</sub>, (x, y)<sub>2</sub>, ..., (x, y)<sub>N</sub> of (discriminating) feature variables x and class labels y
- The learner uses example classes to know how good it is doing

#### **Reinforcement learning**

- Instead of labels, some sort of reward system (e.g. game score)
- Goal: maximise future payoff by optimising decision policy
- May not even "learn" anything from data, but remembers what triggers reward or punishment

#### **Unsupervised** learning

- e.g. clustering: find similarities in training sample, without having predefined categories
- Discover good internal representation of the input
- Not biased by pre-determined classes  $\Rightarrow$  may discover unexpected features!

# Finding the multivariate discriminant y = f(x)



- Given our N examples  $(x, y)_1, \ldots, (x, y)_N$  we need
  - a function class  $\mathbb{F} = \{f(x, w)\}$  (w: parameters of prediction to be found)
  - a constraint Q(w) on  $\mathbb{F}$  (regularisation term)
  - a loss or error function L(y, f), encoding what is lost if f is poorly chosen in  $\mathbb{F}$  (i.e., f(x, w) far from the desired y = f(x))
- Cannot minimise *L* directly (would depend on the dataset used), but rather its average over a training sample, the empirical risk:

$$R(w) = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i, w))$$

subject to constraint Q(w), so we minimise the cost function:

$$C(w) = R(w) + \lambda Q(w)$$

where  $\lambda$  controls the strength of regularisation

• At the minimum of C(w) we select  $f(x, w_*)$ , our estimate of y = f(x)



#### Data generated from an unknown function with unknown noise





#### Constant least squares fit, RMSE = 0.915









Yann Coadou (CPPM) — Machine learning

IDPASC School, online, 7 Sep 2021 11/137















# Representation Class: training



#### Poly(9) least squares fit, RMSE =0





#### Data generated from an unknown function with unknown noise





#### Const. least squares fit, training RMSE = 0.915, test RMSE = 1.067





#### Linear least squares fit, training RMSE = 0.581, test RMSE = 0.734





Quadr. least squares fit, training RMSE = 0.579, test RMSE = 0.723





#### Cubic least squares fit, training RMSE = 0.339, test RMSE = 0.672













## Choice of function class









- Trade-off between approximation error and estimation error
- Take into account sample size
- Measure (and penalise) complexity
- Use independent test sample
- In practice, no need to correctly guess the function class, but need enough flexibility in your model, balanced with complexity cost



### Quadratic and linear discriminants

**(Boosted)** Decision trees

### 6 Neural networks

Deep neural networks



• Suppose densities s(x) and b(x) are multivariate Gaussians:

$$\begin{aligned} \mathsf{Gaussian}(x|\mu,\Sigma) &= \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp^{\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)} \end{aligned}$$
 with vector of means  $\mu$  and covariance matrix  $\Sigma$ 

• Then B(x) = s(x)/b(x) (or its logarithm) can be expressed explicitly:  $\ln B(x) = \lambda(x) \equiv \chi^2(\mu_B, \Sigma_B) - \chi^2(\mu_S, \Sigma_S)$ 



with 
$$\chi^2(\mu, \Sigma) = (x - \mu)^T \Sigma^{-1}(x - \mu)$$

- Fixed value of λ(x) defines quadratic hypersurface partitioning *n*-dimensional space into signal-rich and background-rich regions
- Optimal separation if *s*(*x*) and *b*(*x*) are indeed multivariate Gaussians

# Ruadratic discriminant



'Two moons' data



# 🖉 Quadratic discriminant







# 🖉 Quadratic discriminant







## 🖉 Quadratic discriminant







©Balàzs Kégl



• If in  $\lambda(x)$  the same covariance matrix is used for each class (e.g.  $\Sigma = \Sigma_S + \Sigma_B$ ) one gets Fisher's discriminant:

$$\lambda(x) = w \cdot x$$
 with  $w \propto \Sigma^{-1}(\mu_S - \mu_B)$ 



- Optimal linear separation
- Works only if signal and background have different means!
- Optimal classifier (reaches the Bayes limit) for linearly correlated Gaussian-distributed variables
- Extension to non-linear problems: support vector machines (see • backup)

## 🛹 (Boosted) Decision trees








#### Decision tree origin

- Machine-learning technique, widely used in social sciences. Originally data mining/pattern recognition, then medical diagnosis, insurance/loan screening, etc.

L. Breiman et al., "Classification and Regression Trees" (1984)

### **Basic principle**

- Extend cut-based selection
  - many (most?) events do not have all characteristics of signal or background
  - try not to rule out events failing a particular criterion
- Keep events rejected by one criterion and see whether other criteria could help classify them properly

#### **Binary trees**

- Trees can be built with branches splitting into many sub-branches
- In this lecture: mostly binary trees



#### Start with all events (signal and background) = first (root) node

- sort all events by each variable
- for each variable, find splitting value with best separation between two children
  - mostly signal in one child
  - mostly background in the other
- select variable and splitting value with best separation, produce two branches (nodes)
  - events failing criterion on one side
  - events passing it on the other

### Keep splitting

- Now have two new nodes. Repeat algorithm recursively on each node
- Can reuse the same variable
- Iterate until stopping criterion is reached (min leaf size, max tree depth, insufficient improvement, perfect classification, etc.)
- Splitting stops: terminal node = leaf



 Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>







- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \le p_T^{b_{34}} \le \dots \le p_T^{b_2} \le p_T^{s_{12}}$ •  $H_T^{b_5} \le H_T^{b_3} \le \dots \le H_T^{s_{67}} \le H_T^{s_{43}}$
    - $M_t^{b_6} \le M_t^{s_8} \le \dots \le M_t^{s_{12}} \le M_t^{b_9}$







- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \leq p_T^{b_{34}} \leq \cdots \leq p_T^{b_2} \leq p_T^{s_{12}}$
    - $H_T^{b_5} \leq H_T^{b_3} \leq \cdots \leq H_T^{s_{67}} \leq H_T^{s_{43}}$
    - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7







- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \leq p_T^{b_{34}} \leq \cdots \leq p_T^{b_2} \leq p_T^{s_{12}}$
    - $H_T^{b_5} \leq H_T^{b_3} \leq \cdots \leq H_T^{s_{67}} \leq H_T^{s_{43}}$
    - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7





### Reprision Algorithm example

- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \leq p_T^{b_{34}} \leq \cdots \leq p_T^{b_2} \leq p_T^{s_{12}}$
    - $H_{T_{b_{1}}}^{b_{5}} \leq H_{T_{c_{1}}}^{b_{3}} \leq \cdots \leq H_{T_{c_{1}}}^{s_{67}} \leq H_{T_{c_{1}}}^{s_{43}}$
    - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7
  - split events in two branches: pass or fail  $H_T < 242 \text{ GeV}$





### Reprision Algorithm example

- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \leq p_T^{b_{34}} \leq \cdots \leq p_T^{b_2} \leq p_T^{s_{12}}$
    - $H_T^{b_5} \leq H_T^{b_3} \leq \cdots \leq H_T^{s_{67}} \leq H_T^{s_{43}}$
    - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7
  - split events in two branches: pass or fail  $H_T < 242$  GeV
- Repeat recursively on each node





### Algorithm example

- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:
    - $p_T^{s_1} \leq p_T^{b_{34}} \leq \cdots \leq p_T^{b_2} \leq p_T^{s_{12}}$
    - $H_{T_{f}}^{b_{5}} \leq H_{T}^{b_{3}} \leq \cdots \leq H_{T}^{s_{67}} \leq H_{T_{f}}^{s_{43}}$
    - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7
  - split events in two branches: pass or fail  $H_T < 242 \text{ GeV}$
- Repeat recursively on each node
- Splitting stops: e.g. events with  $H_T < 242$  GeV and  $M_t > 162$  GeV are signal like (p = 0.82)





### Recision tree output

### Run event through tree

- Start from root node
- Apply first best cut
- Go to left or right child node
- Apply best cut for this node
- ...Keep going until...
- Event ends up in leaf

### **DT** Output

р.,<27.6 Fall H,<242

- Purity  $\left(\frac{s}{s+b}\right)$ , with weighted events) of leaf, close to 1 for signal and 0 for background
- or binary answer (discriminant function +1 for signal, -1 or 0 for background) based on purity above/below specified value (e.g.  $\frac{1}{2}$ ) in leaf
- E.g. events with  $H_T < 242$  GeV and  $M_t > 162$  GeV have a DT output of 0.82 or +1

Yann Coadou (CPPM) — Machine learning



### **Optimal split: figure of merit**

- Decrease of impurity for split s of node t into children t<sub>P</sub> and t<sub>F</sub> (goodness of split): Δi(s, t) = i(t) p<sub>P</sub> · i(t<sub>P</sub>) p<sub>F</sub> · i(t<sub>F</sub>)
- Aim: find split  $s^*$  such that  $\Delta i(s^*,t) = \max_{s \in \{\text{splits}\}} \Delta i(s,t)$
- Maximising  $\Delta i(s,t) \equiv$  minimising overall tree impurity



• Also cross section  $\left(-\frac{s^2}{s+b}\right)$  and excess significance  $\left(-\frac{s^2}{b}\right)$ 



# CPPM

#### Reminder

• Need model giving good description of data

#### Reminder

Need model giving good description of data

### Playing with variables

- Number of variables:
  - not affected too much by "curse of dimensionality"
  - CPU consumption scales as *nN* log *N* with *n* variables and *N* training events
- Variable order does not matter: all variables treated equal
- Order of training events is irrelevant (batch training)
- Irrelevant variables:
  - $\bullet\,$  no discriminative power  $\Rightarrow\,$  not used
  - only costs a little CPU time, no added noise
- Can use continuous and discrete variables, simultaneously



### Transforming input variables

- Completely insensitive to replacement of any subset of input variables by (possibly different) arbitrary strictly monotone functions of them (same order ⇒ same DT):
  - $\bullet \ \ \text{convert} \ \ \text{MeV} \to \text{GeV}$
  - no need to make all variables fit in the same range
  - no need to regularise variables (e.g. taking the log)

 $\Rightarrow$  Some immunity against outliers



### Transforming input variables

- Completely insensitive to replacement of any subset of input variables by (possibly different) arbitrary strictly monotone functions of them (same order ⇒ same DT):
  - $\bullet \ \ \text{convert} \ \ \text{MeV} \to \text{GeV}$
  - no need to make all variables fit in the same range
  - no need to regularise variables (e.g. taking the log)

 $\Rightarrow$  Some immunity against outliers

#### Note about actual implementation

- The above is strictly true only if testing all possible cut values
- If there is some computational optimisation (e.g., check only 20 possible cuts on each variable), it may not work anymore



### Variable ranking (mean decrease impurity MDI)

- Ranking of x<sub>i</sub>: add up decrease of impurity each time x<sub>i</sub> is used
- Largest decrease of impurity = best variable

#### Shortcoming: masking of variables

- $x_i$  may be just a little worse than  $x_i$  but will never be picked
- x<sub>j</sub> is ranked as irrelevant
- But remove x<sub>i</sub> and x<sub>i</sub> becomes very relevant
  - $\Rightarrow$  careful with interpreting ranking (specific to training)

#### Permutation importance (mean decrease accuracy MDA)

- Applicable to any already trained classifier
- Randomly shuffle each variable in turn and measure decrease of performance
- Important variable  $\Rightarrow$  big loss of performance
- Can also be performed on validation sample
- Beware of correlations

[Breiman 2001]





### **Choosing variables**

- Usually try to have as few variables as possible
- But difficult: correlations, possibly large number to consider, large phase space with different properties in different regions
- Brute force: with n variables train all n, n-1, etc. combinations, pick best
- Backward elimination: train with n variables, then train all n-1 variables trees and pick best one; now train all n-2 variables trees starting from the n-1 variable list; etc. Pick optimal cost-complexity tree.
- Forward greedy selection: start with k = 1 variable, then train all k + 1 variables trees and pick the best; move to k + 2 variables; etc.







### (Boosted) Decision trees

- Decision trees
- Limitations
  - Training sample composition
  - Ensemble learning
- Boosted decision trees
- Performance examples
- BDTs in real physics cases
- Software and example code



- Small changes in sample can lead to very different tree structures (high variance)
- Performance on testing events may be as good, or not
- Not optimal to understand data from DT rules
- Does not give confidence in result:
  - DT output distribution discrete by nature
  - granularity related to tree complexity
  - tendency to have spikes at certain purity values (or just two delta functions at  $\pm 1$  if not using purity)

## **Tree (in)stability: distributed representation**



- One tree:
  - one information about event (one leaf)
  - cannot really generalise to variations not covered in training set (at most as many leaves as input size)
- Many trees:
  - distributed representation: number of intersections of leaves exponential in number of trees
  - $\bullet\,$  many leaves contain the event  $\Rightarrow$  richer description of input pattern



## Tree (in)stability solution: averaging



### Build several trees and average the output



- K-fold cross-validation (good for small samples)
  - divide training sample  $\mathcal{L}$  in K subsets of equal size:  $\mathcal{L} = \bigcup_{k=1..K} \mathcal{L}_k$
  - Train tree  $T_k$  on  $\mathcal{L} \mathcal{L}_k$ , test on  $\mathcal{L}_k$
  - DT output  $= \frac{1}{K} \sum_{k=1..K} T_k$
- Bagging, boosting, random forests: ensemble learning

### 🖉 Boosted decision trees





### (Boosted) Decision trees

- Decision trees
- Limitations
- Boosted decision trees
  - Introduction
  - AdaBoost
  - Overtraining?
  - Clues to boosting performance
  - Gradient boosting
- Performance examples
- BDTs in real physics cases
- Software and example code

### Roosting: a brief history

### First provable algorithm [Schapire 1990]

- Train classifier  $T_1$  on N events
- Train  $T_2$  on new N-sample, half of which misclassified by  $T_1$
- Build  $T_3$  on events where  $T_1$  and  $T_2$  disagree
- Boosted classifier: MajorityVote(T<sub>1</sub>, T<sub>2</sub>, T<sub>3</sub>)



### Roosting: a brief history

### First provable algorithm [Schapire 1990]

- Train classifier  $T_1$  on N events
- Train  $T_2$  on new N-sample, half of which misclassified by  $T_1$
- Build  $T_3$  on events where  $T_1$  and  $T_2$  disagree
- Boosted classifier: MajorityVote(T<sub>1</sub>, T<sub>2</sub>, T<sub>3</sub>)

#### Then

- Variation [Freund 1995]: boost by majority (combining many learners with fixed error rate)
- Freund&Schapire joined forces: 1<sup>st</sup> functional model AdaBoost (1996)



## Roosting: a brief history

### First provable algorithm [Schapire 1990]

- Train classifier  $T_1$  on N events
- Train  $T_2$  on new N-sample, half of which misclassified by  $T_1$
- Build  $T_3$  on events where  $T_1$  and  $T_2$  disagree
- Boosted classifier: MajorityVote(T<sub>1</sub>, T<sub>2</sub>, T<sub>3</sub>)

#### Then

- Variation [Freund 1995]: boost by majority (combining many learners with fixed error rate)
- Freund&Schapire joined forces: 1<sup>st</sup> functional model AdaBoost (1996)

### When it really picked up in HEP

- MiniBooNe compared performance of different boosting algorithms and neural networks for particle ID [MiniBooNe 2005]
- D0 claimed first evidence for single top quark production [D0 2006]
- CDF copied 😌 (2008). Both used BDT for single top observation







### What is boosting?

- General method, not limited to decision trees
- Hard to make a very good learner, but easy to make simple, error-prone ones (but still better than random guessing)
- Goal: combine such weak classifiers into a new more stable one, with smaller error

### AdaBoost

- Introduced by Freund&Schapire in 1996
- Stands for adaptive boosting
- Learning procedure adjusts to training data to classify it better
- Many variations on the same theme for actual implementation
- Usually leads to better results than without boosting

### AdaBoost algorithm



- Check which events of training sample  $\mathbb{T}_k$  are misclassified by  $T_k$ :
  - $\mathbb{I}(X) = 1$  if X is true, 0 otherwise
  - for DT output in  $\{\pm 1\}$ : isMisclassified<sub>k</sub> $(i) = \mathbb{I}(y_i \times T_k(x_i) \le 0)$
  - or is Misclassified  $_k(i) = \mathbb{I}(y_i \times (T_k(x_i) 0.5) \le 0)$  in purity convention
  - misclassification rate:

$$R(T_k) = \varepsilon_k = \frac{\sum_{i=1}^N w_i^k \times \text{isMisclassified}_k(i)}{\sum_{i=1}^N w_i^k}$$

- Derive tree weight  $\alpha_k = \beta \times \ln((1 \varepsilon_k) / \varepsilon_k)$
- Increase weight of misclassified events in  $\mathbb{T}_k$  to create  $\mathbb{T}_{k+1}$ :

$$w_i^k \to w_i^{k+1} = w_i^k \times e^{\alpha_k}$$

- Train  $T_{k+1}$  on  $\mathbb{T}_{k+1}$
- Boosted result of event *i*:  $T(i) = \frac{1}{\sum_{k=1}^{N_{\text{tree}}} \alpha_k} \sum_{k=1}^{N_{\text{tree}}} \alpha_k T_k(i)$





#### Misclassification rate $\varepsilon$ on training sample



• If each tree has  $\varepsilon_k \neq 0.5$  (i.e. better than random guessing): the error rate falls to zero for sufficiently large  $N_{tree}$ 

• Corollary: training data is overfitted

### **Overtraining**?

- Error rate on test sample may reach a minimum and then potentially rise. Stop boosting at the minimum.
- In principle AdaBoost *must* overfit training sample
- In many cases in literature, no loss of performance due to overtraining
  - may have to do with fact that successive trees get in general smaller and smaller weights
  - trees that lead to overtraining contribute very little to final DT output on validation sample

### Overtraining estimation: good or bad?



Yann Coadou (CPPM) — Machine learning



#### Efficiency vs. background fraction



• Clear overtraining, but still better performance after boosting

## $\gtrsim$ Cross section significance $(s/\sqrt{s+b})$



- More relevant than testing error
- Reaches plateau
- Afterwards, boosting does not hurt (just wasted CPU)
- Applicable to any other figure of merit of interest for your use case



## Relues to boosting performance



- First tree is best, others are minor corrections
- Specialised trees do not perform well on most events  $\Rightarrow$  decreasing tree weight and increasing misclassification rate
- Last tree is not better evolution of first tree, but rather a pretty bad DT that only does a good job on few cases that the other trees could not get right
- But adding trees may increase reliability of prediction: margins explanation [Shapire&Freund 2012]
- Double descent risk curve and interpolation regime [Belkin 2019]

- AdaBoost recast in a statistical framework: corresponds to minimising an exponential loss
- Generalisation: formulate boosting as numerical optimisation problem, minimise loss function by adding trees using gradient descent procedure
- Procedure:
  - Build imperfect model  $F_k$  at step k (sometimes  $F_k(x) \neq y$ )
  - Improve model:  $F_{k+1}(x) = F_k(x) + h_k(x) = y$ , or residual  $h_k(x) = y F_k(x)$
  - Train new classifier on residual
- Example: mean squared error loss function  $L_{MSE}(x, y) = \frac{1}{2} (y - F_k(x))^2$

• minimising loss  $J = \sum_{i} L_{MSE}(x_i, y_i)$  leads to  $\frac{\partial J}{\partial F_k(x_i)} = F_k(x_i) - y_i$ 

 $\Rightarrow$  residual as negative gradient:  $h_k(x_i) = y_i - F_k(x_i) = -\frac{\partial J}{\partial F_k(x_i)}$ 

• Generalised to any differentiable loss function

[Friedman 2001]











- In ROOT using TMVA and create\_circ macro from \$ROOTSYS/tutorials/tmva/createData.C to generate dataset
- Plots: TMVA:::TMVAGui("filename");



### Boosting longer (TMVA: NTrees)

- Compare performance of single DT and BDT with more and more trees (5 to 400)
- All other parameters at TMVA default (would be 400 trees)






- Note: max tree depth = 3
- Single (small) DT: not so good. Note: a larger tree would solve this problem
- More trees ⇒ improve performance (less step-like, closer to optimal separation) until saturation
- Largest BDTs: wiggle a little around the contour
   ⇒ picked up features of training sample, that is, overtraining

## ~ Training/testing output





- Better shape with more trees: quasi-continuous
- Overtraining because of disagreement between training and testing? Let's see...

# Performance in optimal significance





- Best significance actually obtained with last BDT, 400 trees!
- But to be fair, equivalent performance with 10 trees already
- Less "stepped" output desirable?  $\Rightarrow$  maybe 50 is reasonable

# Performance in optimal significance







- Boosting weight decreases fast and stabilises
- First trees have small error fractions, then increases towards 0.5 (random guess)
- ullet  $\Rightarrow$  confirms that best trees are first ones, others are small corrections





#### Bagging (Bootstrap aggregating)

- Before building tree  $T_k$  take random sample of N events from training sample with replacement
- Train  $T_k$  on it
- Events not picked form "out of bag" validation sample
- Applicable to other techniques than DT
  - tends to produce more stable and better classifier



## [Breiman 1996]

#### Bagging (Bootstrap aggregating)

- Before building tree  $T_k$  take random sample of N events from training sample with replacement
- Train  $T_k$  on it
- Events not picked form "out of bag" validation sample
- Applicable to other techniques than DT
  - tends to produce more stable and better classifier

#### Random forests

- Same as bagging
- In addition, pick random subset of variables to consider for each node split
- Two levels of randomisation, much more stable output
- Often as good as boosting

### [Breiman 1996]

#### [Breiman 2001]

## RDTs in real physics cases 🖉





#### Yann Coadou (CPPM) — Machine learning

#### IDPASC School, online, 7 Sep 2021

# BDT in HEP

## ATLAS b-tagging in Run 2

- Run 1 MV1c: NN trained from output of other taggers
- Run 2 MV2c20: BDT using feature variables of underlying algorithms and p<sub>T</sub>, η of jets
- Run 2: introduced IBL (new innermost pixel layer)

 $\Rightarrow$  explains part of the performance gain, but

not all

#### ATLAS $t\bar{t}t\bar{t}$ production evidence

- ► Eur. Phys. J. C 80 (2020) 1085 → arXiv:2007.14858 [hep-ex]
- BDT output used in final fit to measure cross section
- Constraints on systematic uncertainties from profiling







52/137

#### → Eur. Phys. J. C 79 (2019) 970



#### CMS-PAS-HIG-13-001

Hard to use more BDT in an analysis:

- vertex selected with BDT
- 2<sup>nd</sup> vertex BDT to estimate probability to be within 1cm of interaction point
- photon ID with BDT
- photon energy corrected with BDT regression
- event-by-event energy uncertainty from another BDT
- several BDT to extract signal in different categories



# RDT in HEP: final state reconstruction



#### $t\bar{t}H(b\bar{b})$ reconstruction

- Match jets and partons in high-multiplicity final state
- BDT trained on all combinations
- New inputs to classification BDT
- Access to Higgs  $p_T$ , origin of *b*-jets







# BDT and systematics

- No particular rule
- BDT output can be considered as any other cut variable (just more powerful). Evaluate systematics by:
  - varying cut value
  - retraining
  - calibrating, etc.
- Most common (and appropriate): propagate other uncertainties (detector, theory, etc.) up to BDT ouput and check how much the analysis is affected
- More and more common: profiling. Watch out:
  - BDT output powerful
  - signal region (high BDT output) probably low statistics
     ⇒ potential recipe for disaster if modelling is not good
- May require extra systematics, not so much on technique itself, but because it probes specific corners of phase space and/or wider parameter space (usually loosening pre-BDT selection cuts)

# Recision trees are not dead! e.g. NeurIPS2019



- PIDForest: Anomaly Detection via Partial Identification 
   NeurIPS
- A Debiased MDI (Mean Decrease of Impurity) Feature Importance Measure for Random Forests • NeurIPS
- MonoForest framework for tree ensemble analysis
   NeurIPS
- Faster Boosting with Smaller Memory (Yoav S Freund) NeurIPS
- Minimal Variance Sampling in Stochastic Gradient Boosting
   NeurIPS
- Regularized Gradient Boosting NeurIPS
- Partitioning Structure Learning for Segmented Linear Regression Trees • NeurIPS
- Random Tessellation Forests
   NeurIPS
- Optimal Sparse Decision Trees
   NeurIPS
- Provably robust boosted decision stumps and trees against adversarial attacks • NeurIPS
- Robustness Verification of Tree-based Models 

   NeurIPS

- Go for a fully integrated solution
  - use different multivariate techniques easily
  - spend your time on understanding your data and model
- Examples:
  - TMVA (Toolkit for MultiVariate Analysis) Integrated in ROOT, complete manual
    - Example code in backup
  - scikit-learn (python)

#### • Dedicated to BDT but transparently integrated with e.g. scikit-learn:

- XGBoost (popular in HEP) arXiv:1603.02754 (note: cannot handle negative weights)
- LightGBM (Microsoft)
- CatBoost (Yandex)
- Several examples in IDPASC ML hands-on session

https://github.com/dmlc/xgboost





L. Breiman, J.H. Friedman, R.A. Olshen and C.J. Stone, Classification and Regression Trees, Wadsworth, Stamford, 1984 R.E. Schapire, "The strength of weak learnability" Machine Learning 5 (1990) 197 Y. Freund, "Boosting a weak learning algorithm by majority" Information and computation 121 (1995) 256 Y. Freund and R.E. Schapire, "Experiments with a New Boosting Algorithm" in Machine Learning: Proceedings of the Thirteenth International Conference, edited by L. Saitta (Morgan Kaufmann, San Fransisco, 1996) p. 148 Y. Freund and R.E. Schapire, "A short introduction to boosting" R. E. Schapire and Y. Freund, "Boosting: Foundations and Algorithms", MIT Press, 2012. Y. Freund and R.E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting" ( Journal of Computer and System Sciences 55 (1997) 119 J.H. Friedman, T. Hastie and R. Tibshirani, "Additive logistic regression: a statistical view of boosting" Annals of Statistics 28 (2000) 377









# Backup

Yann Coadou (CPPM) — Machine learning

IDPASC School, online, 7 Sep 2021 138/137

# Support vector machines 🖉

- CPPM
- Fisher discriminant: may fail completely for highly non-Gaussian densities
- But linearity is good feature  $\Rightarrow$  try to keep it
- Generalising Fisher discriminant: data non-separable in *n*-dim space  $\mathbb{R}^n$ , but better separated if mapped to higher dimension space  $\mathbb{R}^H$ :  $h: x \in \mathbb{R}^n \to z \in \mathbb{R}^H$
- Use hyper-planes to partition higher dim space:  $f(x) = w \cdot h(x) + b$
- Example:  $h: (x_1, x_2) \rightarrow (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$







• Consider separable data in  $\mathbb{R}^{H}$ , and three parallel hyper-planes:  $w \cdot h(x) + b = 0$  (separating hyper-plane between red and blue)  $w \cdot h(x_1) + b = +1$  (contains  $h(x_1)$ )

 $w \cdot h(x_2) + b = -1$  (contains  $h(x_2)$ )



- Subtract blue from red:  $w \cdot (h(x_1) - h(x_2)) = 2$
- With unit vector  $\hat{w} = w/||w||$ :  $\hat{w} \cdot (h(x_1) - h(x_2)) = 2/||w|| = m$
- Margin *m* is distance between red and blue planes
- Best separation: maximise margin
- $\Rightarrow$  empirical risk margin to minimise:  $R(w) \propto ||w||^2$



- When minimising R(w), need to keep signal and background separated
- Label red dots y = +1 ("above" red plane) and blue dots y = -1 ("below" blue plane)

# • Since: $w \cdot h(x) + b > 1$ for red dots $w \cdot h(x) + b < -1$ for blue dots

all correctly classified points will satisfy constraints:

$$y_i(w \cdot h(x_i) + b) \geq 1, \ \forall i = 1, \dots, N$$

• Using Lagrange multipliers  $\alpha_i > 0$ , cost function can be written:

$$C(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^{N} \alpha_i \left[ y_i \left( w \cdot h(x_i) + b \right) - 1 \right]$$



#### Minimisation

• Minimise cost function  $C(w, b, \alpha)$  with respect to w and b:

$$C(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j (h(x_i) \cdot h(x_j))$$

 At minimum of C(α), only non-zero α<sub>i</sub> correspond to points on red and blue planes: support vectors

#### **Kernel functions**

Issues:

- need to find *h* mappings (potentially of infinite dimension)
- need to compute scalar products  $h(x_i) \cdot h(x_j)$
- Fortunately  $h(x_i) \cdot h(x_j)$  are equivalent to some kernel function  $K(x_i, x_j)$  that does the mapping and the scalar product:

$$C(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$





- In reality: do not know a priori the right kernel
- $\bullet\,\Rightarrow\,$  have to test different standard kernels and use the best one

Support vector machines: non-separable data 🖉



- Even in infinite dimension space, data are often non-separable
- Need to relax constraints:

$$y_i(w \cdot h(x_i) + b) \geq 1 - \xi_i$$



with slack variables  $\xi_i > 0$ 

- C(w, b, α, ξ) depends on ξ, modified C(α, ξ) as well
- Values determined during minimisation



#### Defined for many classes

• Gini = 
$$\sum_{i,j\in\{\text{classes}\}}^{i\neq j} p_i p_j$$

#### Statistical interpretation

- Assign random object to class i with probability  $p_i$ .
- Probability that it is actually in class j is p<sub>j</sub>
- $\bullet \Rightarrow \mathsf{Gini} = \mathsf{probability} \text{ of misclassification}$

#### For two classes (signal and background)

• 
$$i = s, b$$
 and  $p_s = p = 1 - p_b$ 

• 
$$\Rightarrow$$
 Gini =  $1 - \sum_{i=s,b} p_i^2 = 2p(1-p) = \frac{2sb}{(s+b)^2}$ 

- Most popular in DT implementations
- Usually similar performance to e.g. entropy

# Runing a tree I

#### Why prune a tree?

- Possible to get a perfect classifier on training events
- Mathematically misclassification error can be made as little as wanted
- E.g. tree with one class only per leaf (down to 1 event per leaf if necessary)
- Training error is zero
- But run new independent events through tree (testing or validation sample): misclassification is probably > 0, overtraining
- Pruning: eliminate subtrees (branches) that seem too specific to training sample:
  - a node and all its descendants turn into a leaf

#### Pruning algorithms

- Pre-pruning (early stopping condition like min leaf size, max depth)
- Expected error pruning (based on statistical error estimate)
- Cost-complexity pruning (penalise "complex" trees with many nodes/leaves)



#### **Pre-pruning**

- Stop tree growth during building phase
- Already seen: minimum leaf size, minimum separation improvement, maximum depth, etc.
- Careful: early stopping condition may prevent from discovering further useful splitting

#### Expected error pruning

- Grow full tree
- When result from children not significantly different from result of parent, prune children
- Can measure statistical error estimate with binomial error  $\sqrt{p(1-p)/N}$  for node with purity p and N training events
- No need for testing sample
- Known to be "too aggressive"





- Idea: penalise "complex" trees (many nodes/leaves) and find compromise between good fit to training data (larger tree) and good generalisation properties (smaller tree)
- With misclassification rate R(T) of subtree T (with N<sub>T</sub> nodes) of fully grown tree T<sub>max</sub>:

cost complexity  $R_{\alpha}(T) = R(T) + \alpha N_T$ 

- $\alpha = \text{ complexity parameter}$
- Minimise  $R_{\alpha}(T)$ :
  - small  $\alpha$ : pick  $T_{max}$
  - large  $\alpha :$  keep root node only,  $\mathcal{T}_{\textit{max}}$  fully pruned
- First-pass pruning, for terminal nodes  $t_L$ ,  $t_R$  from split of t:
  - by construction  $R(t) \ge R(t_L) + R(t_R)$
  - if  $R(t) = R(t_L) + R(t_R)$  prune off  $t_L$  and  $t_R$

Pruning a tree IV: cost-complexity pruning



- For node t and subtree T<sub>t</sub>:
  - if t non-terminal,  $R(t) > R(T_t)$  by construction
  - $R_{\alpha}(\{t\}) = R_{\alpha}(t) = R(t) + \alpha \ (N_{T} = 1)$
  - if  $R_{\alpha}(T_t) < R_{\alpha}(t)$  then branch has smaller cost-complexity than single node and should be kept
  - at critical  $\alpha = \rho_t$ , node is preferable
  - to find  $\rho_t$ , solve  $R_{\rho_t}(T_t) = R_{\rho_t}(t)$ , or:  $\rho_t = \frac{R(t) R(T_t)}{N_{\tau} 1}$
  - node with smallest  $\rho_t$  is weakest link and gets pruned
  - apply recursively till you get to the root node
- This generates sequence of decreasing cost-complexity subtrees
- Compute their true misclassification rate on validation sample:
  - will first decrease with cost-complexity
  - then goes through a minimum and increases again
  - pick this tree at the minimum as the best pruned tree
- Note: best pruned tree may not be optimal in a forest

# Introduction to TMVA

• TMVA: Toolkit for MultiVariate Analysis

https://root.cern/tmva

https://github.com/root-project/root/tree/master/tmva

(ROOT v6.20.06)

- Written by physicists
- In C++ (also python API), integrated in ROOT
- Quite complete manual
- Includes many different multivariate/machine learning techniques
- To compile, add appropriate header files in your code (e.g., #include "TMVA/Factory.h") and this to your compiler command line: 'root-config --cflags --libs' -1TMVA
- More complete examples of code: \$ROOTSYS/tutorials/tmva
  - createData.C macro to make example datasets
  - classification and regression macros
  - also includes Keras examples (deep learning)
- Sometimes useful performance measures (more in these headers): #include "TMVA/ROCCalc.h" TMVA::ROCCalc(TH1\* S,TH1\* B).GetROCIntegral();

#include "TMVA/Tools.h"

TMVA::gTools().GetSeparation(TH1\* S,TH1\* B);





TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile, "!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");



TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root") 
Trree* sig = (TTree*)inputFile->Get("TreeS");
Trree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
    new TMVA::DataLoader("dataset"):
```

dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight);





TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root") 
Trree* sig = (TTree*)inputFile->Get("TreeS");
Trree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
    new TMVA::DataLoader("dataset");
dataloader->AddSignalTree(sig, sigWeight);
dataloader->AddBackgroundTree(bkg, bkgWeight);
dataloader->AddVariable("var0", 'F');
dataloader->AddVariable("var1", 'F');
```



TCut mycut = "";



TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root") TTree* sig = (TTree*)inputFile->Get("TreeS");
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
```

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = "";



dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random");



TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root") -
TTree* sig = (TTree*)inputFile->Get("TreeS");
                                                       3
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
                                                       0F
   new TMVA::DataLoader("dataset");
dataloader->AddSignalTree(sig, sigWeight);
                                                       -2
dataloader->AddBackgroundTree(bkg, bkgWeight);
dataloader->AddVariable("var0", 'F');
dataloader->AddVariable("var1", 'F');
TCut mycut = "";
dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random");
factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:
   MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80");
factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher");
```

TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisTvpe=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root") 5
TTree* sig = (TTree*)inputFile->Get("TreeS");
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
```

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = "";



dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random"); factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:

MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80");
factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher");
factory->TrainAllMethods(); // Train MVAs using training events

TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisTvpe=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root")
TTree* sig = (TTree*)inputFile->Get("TreeS");
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
```

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = ""; dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random");



dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random"); factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:

MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80");
factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher");
factory->TrainAllMethods(); // Train MVAs using training events
factory->TestAllMethods(); // Evaluate all MVAs using test events
TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification");

```
TFile* inputFile = new TFile("dataSchachbrett.root")
TTree* sig = (TTree*)inputFile->Get("TreeS");
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
TMVA::DataLoader *dataloader =
```

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = "";



dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random"); factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:

MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80");
factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher");
factory->TrainAllMethods(); // Train MVAs using training events
factory->TestAllMethods(); // Evaluate all MVAs using test events
// ----- Evaluate and compare performance of all configured MVAs
factory->EvaluateAllMethods();

TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification"); TFile\* inputFile = new TFile("dataSchachbrett.root"); TTree\* sig = (TTree\*)inputFile->Get("TreeS"); TTree\* bkg = (TTree\*)inputFile->Get("TreeB"); double sigWeight = 1.0; double bkgWeight = 1.0;

TMVA::DataLoader \*dataloader =

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = "";



dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random"); factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:

MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80");
factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher");
factory->TrainAllMethods(); // Train MVAs using training events
factory->TestAllMethods(); // Evaluate all MVAs using test events
// ----- Evaluate and compare performance of all configured MVAs
factory->EvaluateAllMethods();

auto c1 = factory->GetROCCurve(dataloader); // Eager to compare performance

TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification"); TFile\* inputFile = new TFile("dataSchachbrett.root"); TTree\* sig = (TTree\*)inputFile->Get("TreeS"); TTree\* bkg = (TTree\*)inputFile->Get("TreeB"); double sigWeight = 1.0; double bkgWeight = 1.0; TMVA::DataLoader \*dataloader =

new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F'); dataloader->AddVariable("var1", 'F'); TCut mycut = ""; dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random");



factory->BookMethod(dataloader, TMVA::Types::kFisher", "Fisher", "!H:!V:Fisher"); factory->TrainAllMethods(); // Train MVAs using training events factory->TestAllMethods(); // Evaluate all MVAs using test events // ----- Evaluate and compare performance of all configured MVAs factory->EvaluateAllMethods(); auto c1 = factory->GetROCCurve(dataloader); // Eager to compare performance outputFile->Close();

delete factory; delete dataloader;

TFile\* outputFile = TFile::Open("output.root", "RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

3

0

"!V:Color:DrawProgressBar:Transformations=I:AnalysisType=Classification"); TFile\* inputFile = new TFile("dataSchachbrett.root") -TTree\* sig = (TTree\*)inputFile->Get("TreeS"); TTree\* bkg = (TTree\*)inputFile->Get("TreeB"); double sigWeight = 1.0; double bkgWeight = 1.0; TMVA::DataLoader \*dataloader =

```
new TMVA::DataLoader("dataset");
                                                        -1F
dataloader->AddSignalTree(sig, sigWeight);
                                                       -2F
dataloader->AddBackgroundTree(bkg, bkgWeight);
dataloader->AddVariable("var0", 'F');
dataloader->AddVariable("var1", 'F');
TCut mycut = "";
dataloader->PrepareTrainingAndTestTree(mycut,"SplitMode=Random");
factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400:
```

MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80"); factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher"); factory->TrainAllMethods(); // Train MVAs using training events factory->TestAllMethods(); // Evaluate all MVAs using test events // ---- Evaluate and compare performance of all configured MVAs factory->EvaluateAllMethods(); auto c1 = factory->GetROCCurve(dataloader); // Eager to compare performance outputFile->Close(); delete factory; delete dataloader; TMVA::TMVAGui("output.root");



TFile\* inputFile = new TFile("dataSchachbrett.root"); TTree\* data = (TTree\*)inputFile->Get("TreeS"); Float\_t var0=-99., var1=-99.; data->SetBranchAddress("var0", &var0); data->SetBranchAddress("var1", &var1);



```
TFile* inputFile = new TFile("dataSchachbrett.root");
TTree* data = (TTree*)inputFile->Get("TreeS");
Float_t var0=-99., var1=-99.;
data->SetBranchAddress("var0", &var0);
data->SetBranchAddress("var1", &var1);
TMVA::Reader *reader = new TMVA::Reader();
reader->AddVariable( "var0", &var0 );
reader->AddVariable( "var1", &var1 );
```



```
TFile* inputFile = new TFile("dataSchachbrett.root");
TTree* data = (TTree*)inputFile->Get("TreeS");
Float_t var0=-99., var1=-99.;
data->SetBranchAddress("var0", &var0);
data->SetBranchAddress("var1", &var1);
TMVA::Reader *reader = new TMVA::Reader();
reader->AddVariable( "var0", &var0 );
reader->AddVariable( "var1", &var1 );
reader->BookMVA( "My BDT", "dataset/weights/TMVAClassification_BDT.weights.xml");
reader->BookMVA( "Fisher discriminant",
    "dataset/weights/TMVAClassification_Fisher.weights.xml");
```

```
TFile* inputFile = new TFile("dataSchachbrett.root");
TTree* data = (TTree*)inputFile->Get("TreeS");
Float t var0=-99.. var1=-99.:
data->SetBranchAddress("var0", &var0);
data->SetBranchAddress("var1", &var1);
TMVA::Reader *reader = new TMVA::Reader();
reader->AddVariable( "var0", &var0 );
reader->AddVariable( "var1", &var1 );
reader->BookMVA( "My BDT", "dataset/weights/TMVAClassification_BDT.weights.xml");
reader->BookMVA( "Fisher discriminant".
  "dataset/weights/TMVAClassification_Fisher.weights.xml");
// ----- start your event loop
for (Long64_t ievt=0; ievt<10; ++ievt) {</pre>
  data->GetEntry(ievt);
  double bdt = reader->EvaluateMVA("My BDT");
  double fisher = reader->EvaluateMVA("Fisher discriminant");
  cout<<"var0="<<var0<<" var1="<<var1<<" BDT="<<bdt<<" Fisher="<<fisher<<end1:</pre>
7
delete reader:
inputFile->Close();
```

```
TFile* inputFile = new TFile("dataSchachbrett.root");
TTree* data = (TTree*)inputFile->Get("TreeS");
Float_t var0=-99., var1=-99.;
data->SetBranchAddress("var0", &var0);
data->SetBranchAddress("var1", &var1);
TMVA::Reader *reader = new TMVA::Reader();
reader->AddVariable( "var0", &var0 );
reader->AddVariable( "var1", &var1 );
reader->BookMVA( "My BDT", "dataset/weights/TMVAClassification_BDT.weights.xml");
reader->BookMVA( "Fisher discriminant".
  "dataset/weights/TMVAClassification_Fisher.weights.xml");
// ----- start your event loop
for (Long64_t ievt=0; ievt<10; ++ievt) {</pre>
 data->GetEntry(ievt);
  double bdt = reader->EvaluateMVA("My BDT");
  double fisher = reader->EvaluateMVA("Fisher discriminant"):
 cout<<"var0="<<var0<<" yar1="<<var1<<" BDT="<<bdt<<" Fisher="<<fisher<<end1:
3
delete reader:
inputFile->Close();
```

• More complete tutorials:

https://github.com/Imoneta/tmva-tutorial

## $\geq$ Compiling TMVA with C++

- To make code compilable (and MUCH faster)
  - Need ROOT and TMVA corresponding header files
  - e.g., for Train.C:

• Compilation:

g++ Train.C 'root-config --cflags --libs' -lTMVA -lTMVAGui -o TMVATrainer

- Train.C: file to compile
- TMVATrainer: name of executable
- -ITMVAGui: just because of TMVA:::TMVAGui("output.root");

#### TMVA: training refinements 🖉

- Common technique: train on even event numbers, test on odd event numbers (and vice versa)
- Can also think of more than two-fold
- Achieve in TMVA by replacing:

```
dataloader->AddSignalTree(sig, sigWeight);
dataloader->AddBackgroundTree(bkg, bkgWeight);
```

• with:

```
TString trainString = "(eventNumber % 2 == 0)";
TString testString = "!"+trainString;
dataloader->AddTree(sig, "Signal", sigWeight, trainString.Data(), "Training");
dataloader->AddTree(sig, "Signal", sigWeight, testString.Data(), "Test");
dataloader->AddTree(bkg, "Background", bkgWeight, trainString.Data(), "Training");
dataloader->AddTree(bkg, "Background", bkgWeight, testString.Data(), "Test");
```

• Use individual event weights:

string eventWeight = "TMath::Abs(eventWeight)"; //Compute event weight
dataloader->SetSignalWeightExpression(eventWeight);
dataloader->SetBackgroundWeightExpression(eventWeight); //Can differ