# Machine learning

Yann Coadou

**CPPM** Marseille



European School of Instrumentation in Particle & Astroparticle Physics

Online edition, 26 January 2022







# Outline

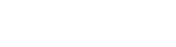
lbab

APPLY NOW





- 2 Optimal discrimination
- **3** Machine learning algorithms
  - Random grid search
  - Genetic algorithms
  - Quadratic and linear discriminants
  - Kernel density estimation
  - Decision trees
  - Boosted decision trees
  - Neural networks
  - Deep neural networks
- 4 Machine learning and (particle) physics
- **5** Machine learning, physics and society
- 6 Conclusion
- 7 Software
- 8 References



**Online format** 

INSTRUMENTATION IN PARTICLE &

### **Typical problems**

- Classification of objects
   separate real and fake leptons/jets/etc.
- Signal enhancement relative to background
  - e.g. galaxy deblending
- Regression: best estimation of a parameter
  - lepton energy, E<sup>miss</sup> value, invariant mass, proportions of radioelements, 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, ...)



# Introduction



#### Input information from various sources

- Kinematic variables (masses, momenta, decay angles, ...)
- Event properties (jet multiplicity, sum of charges, sphericity, ...)
- Detector response (silicon hits, dE/dx, Cherenkov angle, shower profiles, muon hits, ...)
- Telescope images in visible light/infrared/X-ray/etc.

### 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  $f(x) \le q$



 "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E." Tom Mitchell (1997)

 "Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed"

# Machine learning: (un)supervised learning



#### Supervised learning

- Labelled training events: N examples  $(x, y)_1, (x, y)_2, ..., (x, y)_N$  of (discriminating) feature variables x and class labels y
- Use examples' classes to know how good algorithm is doing

# Machine learning: (un)supervised learning



#### Supervised learning

- Labelled training events: N examples  $(x, y)_1, (x, y)_2, \dots, (x, y)_N$  of (discriminating) feature variables x and class labels y
- Use examples' classes to know how good algorithm 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: learn policy

# Machine learning: (un)supervised learning



### Supervised learning

- Labelled training events: 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
- Use examples' classes to know how good algorithm is doing

#### **Reinforcement learning**

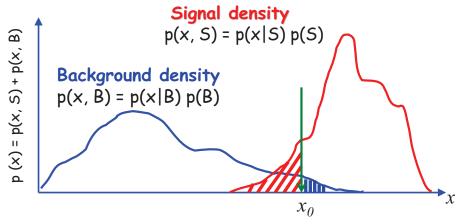
- Instead of labels, some sort of reward system (e.g. game score)
- Goal: maximise future payoff by optimising decision policy
- May not even "learn" anything from data, but remembers what triggers reward or punishment: learn policy

#### **Unsupervised** learning

- Find similarities in training sample, without having predefined categories (no labels/targets)
- Discover good internal representation of the input
- Not biased by pre-determined classes ⇒ may discover unexpected features!



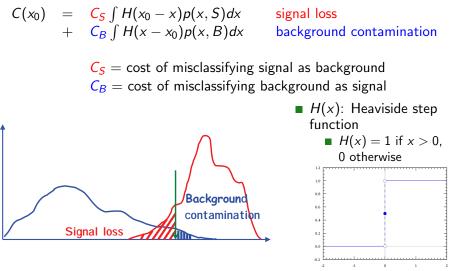
• Where to place a cut  $x_0$  on variable x?



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

# Optimal discrimination: cost of misclassification





Optimal choice: when cost function *C* is minimum

# **Optimal discrimination: Bayes discriminant**



#### Minimising the cost

Minimise

 $C(x_0) = C_S \int H(x_0 - x)p(x, S)dx + C_B \int H(x - x_0)p(x, B)dx$ with respect to the boundary  $x_0$ :

$$0 = C_S \int \delta(x_0 - x) p(x, S) dx - C_B \int \delta(x - x_0) p(x, B) dx$$
  
=  $C_S p(x_0, S) - C_B p(x_0, B)$ 

This gives the Bayes discriminant:

$$BD = \frac{C_B}{C_S} = \frac{p(x_0, S)}{p(x_0, B)} = \frac{p(x_0|S)p(S)}{p(x_0|B)p(B)}$$

# **Optimal 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)

# Machine learning: learning from examples

# СРРМ

### 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, redshift 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
    - $\Rightarrow {\it generalisation \ ability}$

# 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:

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

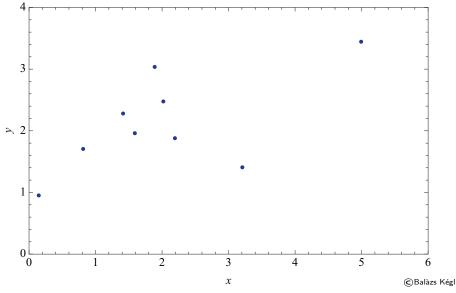
subject to constraint Q(w), so we minimise the cost function:

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

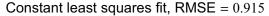
where  $\lambda$  controls the strength of regularisation • At the minimum of C(w) we select  $f(x, w_*)$ , our estimate of y = f(x)

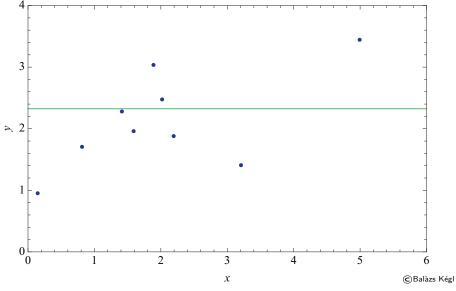


### Data generated from an unknown function with unknown noise

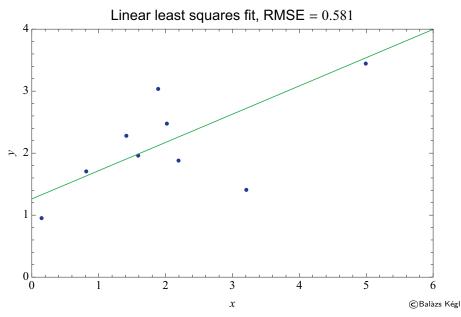




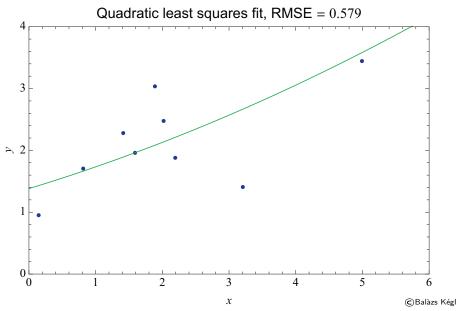




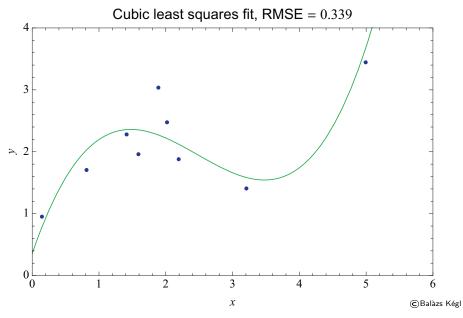




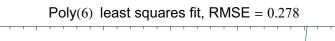


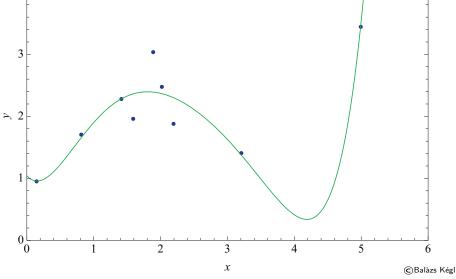




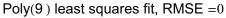


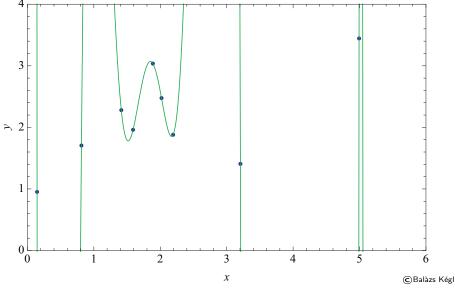










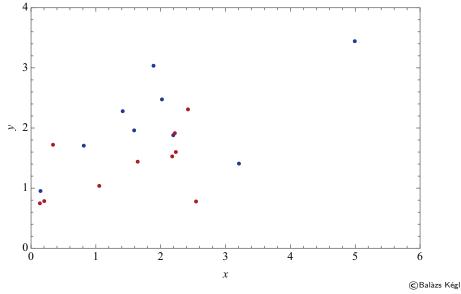




### What happens if running on new data points?

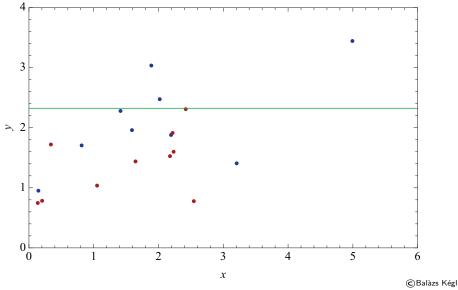


#### Data generated from an unknown function with unknown noise



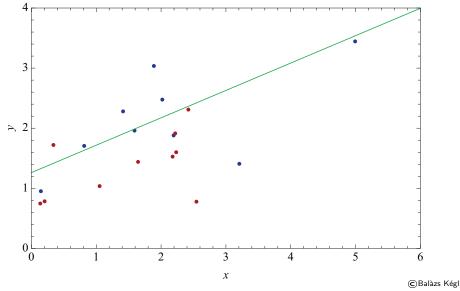


Const. least squares fit, training RMSE = 0.915, test RMSE = 1.067



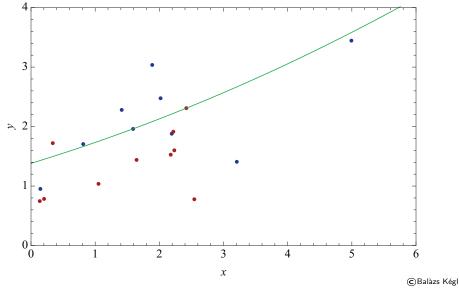


Linear least squares fit, training RMSE = 0.581, test RMSE = 0.734



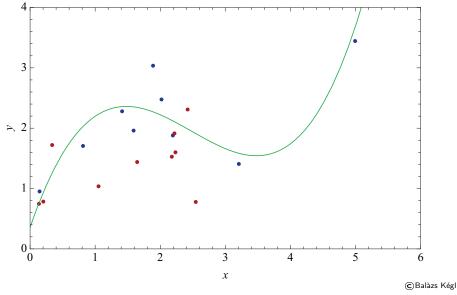


Quadr. least squares fit, training RMSE = 0.579, test RMSE = 0.723



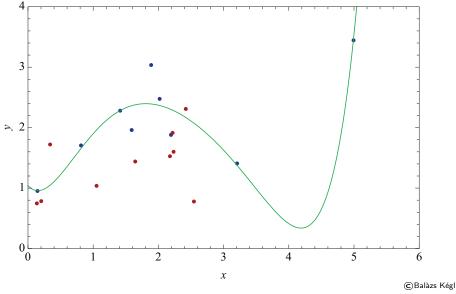


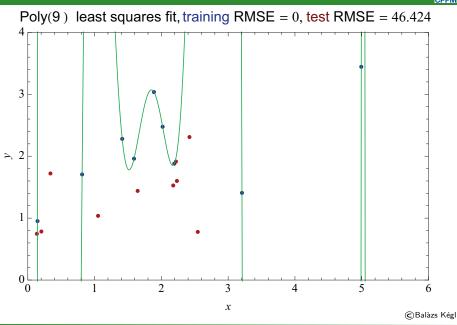
Cubic least squares fit, training RMSE = 0.339, test RMSE = 0.672





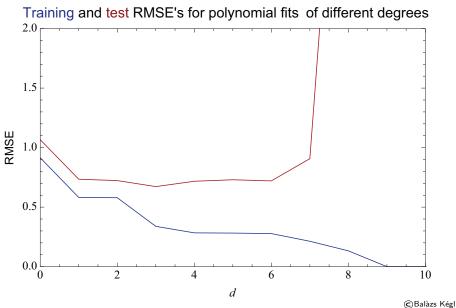






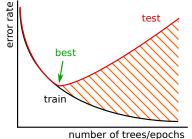
### **Choice of function class**





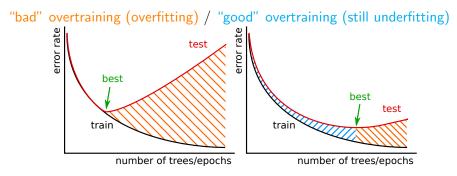


"bad" overtraining (overfitting) / "good" overtraining (still underfitting)

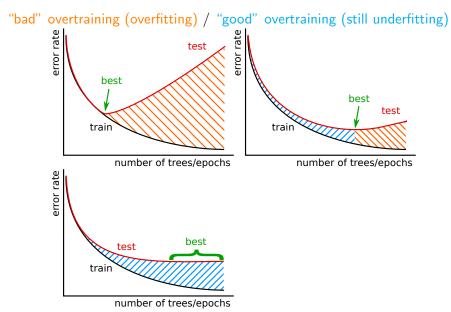


Yann Coadou (CPPM) — Machine learning

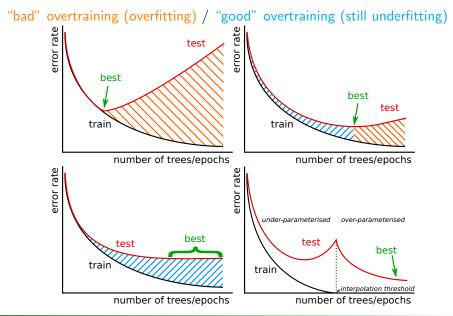










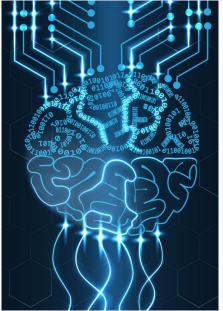




- Trade-off between approximation error and estimation error
- Take into account sample size
- Measure (and penalise) complexity
- Use independent test sample
- In practice, no need to correctly guess the function class, but need enough flexibility in your model, balanced with complexity cost

## Machine learning algorithms





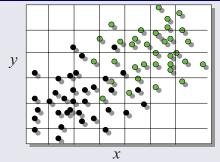
- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Kernel density estimation
- Decision trees
- Boosted decision trees
- Neural networks
- Deep neural networks

### Cut-based analysis and grid search

#### **Cut-based analysis**

- Simple approach: cut on each discriminating variable
- Difficulty: how to optimise the cuts?

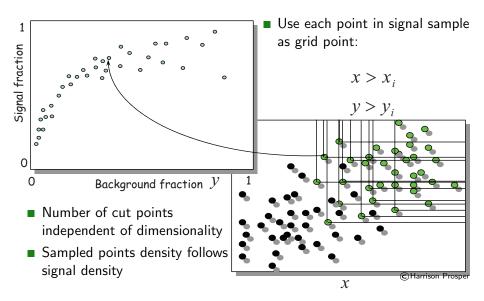
#### Grid search



- Split each variable in K values
- Apply cuts at each grid point: x > x<sub>i</sub>, y > y<sub>i</sub>
- Number of points scales with K<sup>n</sup>: curse of dimensionality

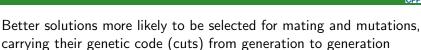




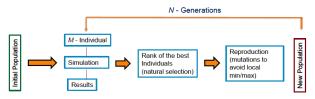




- Inspired by biological evolution
- Model: group (population) of abstract representations (genome/discriminating variables) of possible solutions (individuals/list of cuts)
- Typical processes at work in evolutionary processes:
  - inheritance
  - mutation
  - sexual recombination (a.k.a. crossover)
- Fitness function: value representing the individual's goodness, or comparison of two individuals
- For cut optimisation:
  - good background rejection and high signal efficiency
  - compare individuals in each signal efficiency bin and keep those with higher background rejection



- Algorithm:
  - **1** Create initial random population (cut ensemble)
  - 2 Select fittest individuals
  - **3** Create offsprings through crossover (mix best cuts)
  - 4 Mutate randomly (change some cuts of some individuals)
  - **5** Repeat from 2 until convergence (or fixed number of generations)



■ Good fitness at one generation ⇒ average fitness in the next
 ■ Algorithm focuses on region with higher potential improvement

## Quadratic discriminants: Gaussian problem

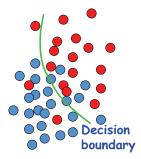


• Suppose densities s(x) and b(x) are multivariate Gaussians:

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

with vector of means  $\mu$  and covariance matrix  $\Sigma$ 

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

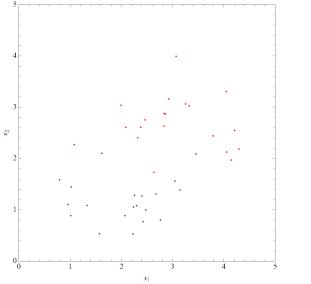


with 
$$\chi^2(\mu, \Sigma) = (x - \mu)^T \Sigma^{-1}(x - \mu)$$

- 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



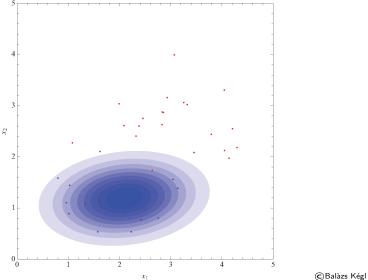
'Two moons' data



©Balàzs Kégl

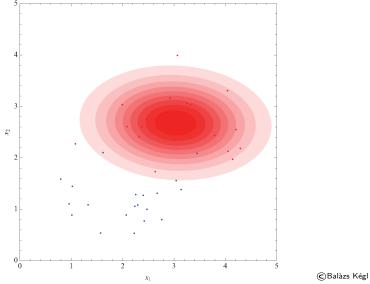




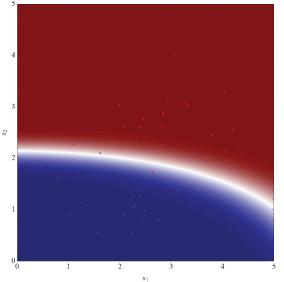












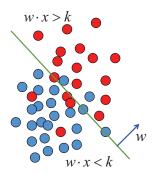
Discriminant function with Gaussian fits

©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  $w \propto \Sigma^{-1}(\mu_S - \mu_B)$ 

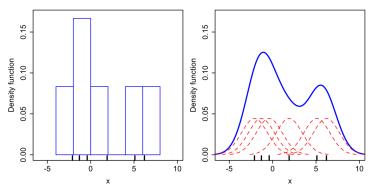


- Optimal linear separation
- Works only if signal and background have different means!
- Optimal classifier (reaches the Bayes limit) for linearly correlated Gaussian-distributed variables
- Extension to non-linear problems: support vector machines (see • backup)



- Introduced by E. Parzen in the 1960s
- Place a kernel  $K(x, \mu)$  at each training point  $\mu$
- Density p(x) at point x approximated by:

$$p(x) \approx \hat{p}(x) = \frac{1}{N} \sum_{j=1}^{N} K(x, \mu_j)$$



#### Choice of kernel

- Any kernel can be used
- In practice, often product of Gaussians:

$$K(x,\mu) = \prod_{i=1}^{n} Gaussian(x_i|\mu, h_i)$$

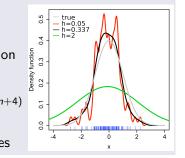
each with bandwidth (width)  $h_i$ 

#### **Optimal bandwidth**

- Too narrow: noisy approximation
- Too wide: loose fine structure
- In principle found by minimising risk function  $R(\hat{p}, p) = \int (\hat{p}(x) p(x))^2 dx$
- For Gaussian densities:

$$h = \sigma \left(\frac{4}{(n+2)N}\right)^{1/(n+2)}$$

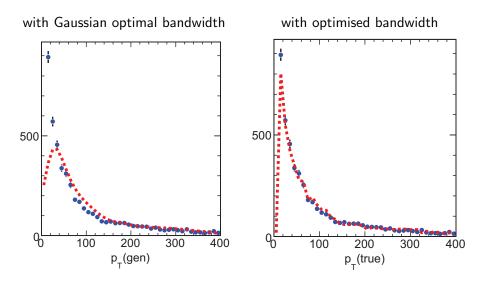
Far from optimal for non-Gaussian densities





## Kernel density estimation (KDE): example





# CPPM

#### Why does it work?

• When  $N \to \infty$ :

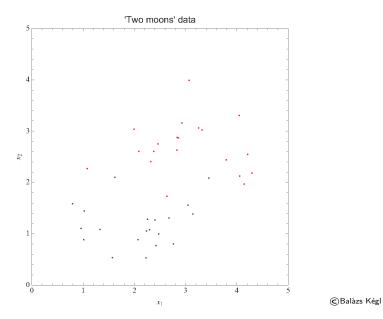
$$\hat{p}(x) = \int K(x,\mu) p(\mu) d\mu$$

- $p(\mu)$ : true density of x
- Kernel bandwidth getting smaller with N, so when  $N \to \infty$ ,  $K(x, \mu) \to \delta^n(x \mu)$  and  $\hat{p}(x) = p(x)$
- KDE gives consistent estimate of probability density p(x)

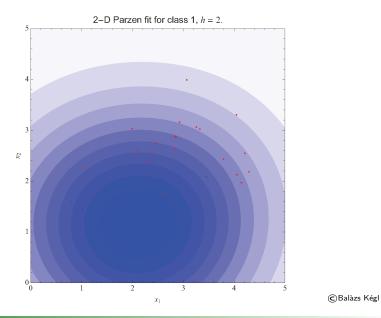
#### Limitations

- Choice of bandwidth non-trivial
- Difficult to model sharp structures (e.g. boundaries)
- Kernels too far apart in regions of low point density
- (both can be mitigated with adaptive bandwidth choice)
- Requires evaluation of N n-dimensional kernels

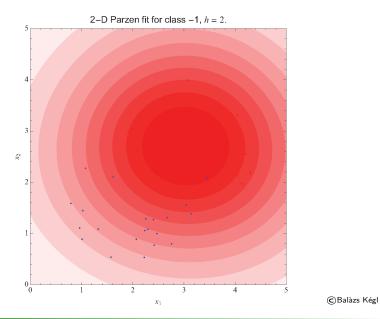




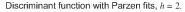


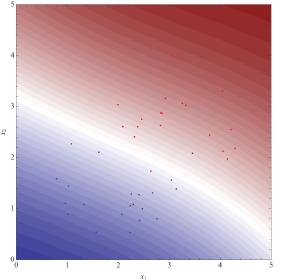








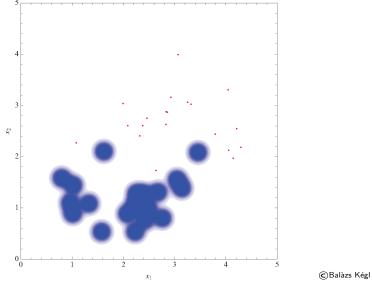




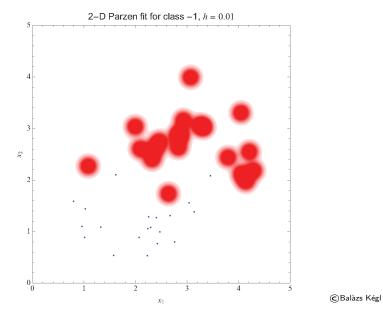
©Balàzs Kégl





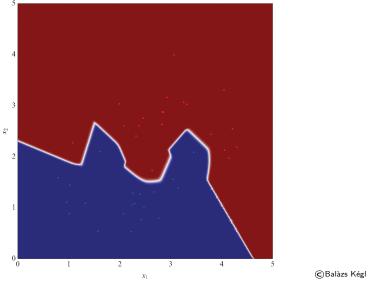






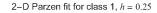


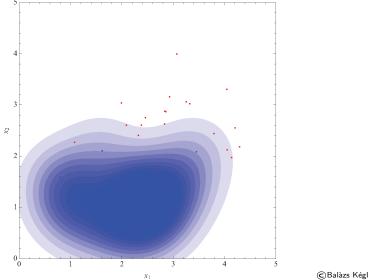
Discriminant function with Parzen fits, h = 0.01



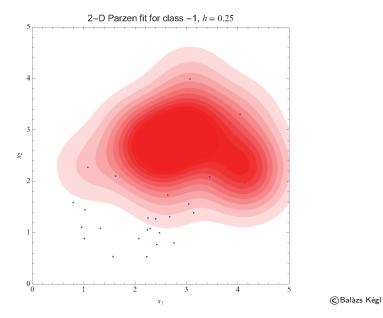
Yann Coadou (CPPM) — Machine learning





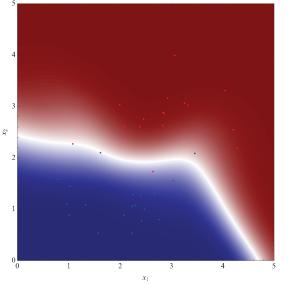








Discriminant function with Parzen fits, h = 0.25

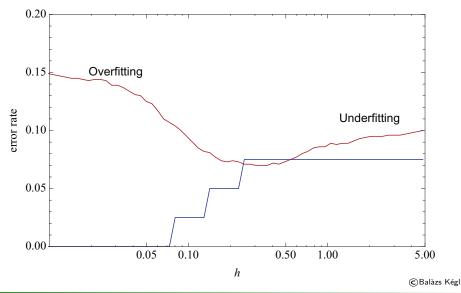


©Balàzs Kégl

#### KDE: choice of bandwidth

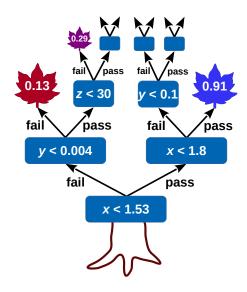


#### Training and test error rates



#### **Decision trees**





- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Kernel density estimation
- Decision trees
  - Algorithm
  - Splitting a node
  - Variable selection
  - Limitations
- Boosted decision trees
- Neural networks
- Deep neural networks

## Introduction



#### 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

Yann Coadou (CPPM) — Machine learning

## Tree building algorithm

#### 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:





- 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:
  - $p_{T}^{s_{1}} \leq p_{T}^{b_{24}} \leq \cdots \leq p_{T}^{b_{2}} \leq p_{T}^{s_{12}}$   $H_{T}^{b_{5}} \leq H_{T}^{b_{3}} \leq \cdots \leq H_{T}^{s_{67}} \leq H_{T}^{s_{43}}$   $M_{t}^{b_{6}} \leq M_{t}^{s_{8}} \leq \cdots \leq M_{t}^{s_{12}} \leq M_{t}^{b_{9}}$  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:
    - $p_{\mathsf{T}}^{s_1} \le p_{\mathsf{T}}^{b_{24}} \le \dots \le p_{\mathsf{T}}^{b_2} \le p_{\mathsf{T}}^{s_{12}}$  $H_{\mathsf{T}}^{b_5} \le H_{\mathsf{T}}^{b_3} \le \dots \le H_{\mathsf{T}}^{s_{67}} \le H_{\mathsf{T}}^{s_{43}}$  $M_t^{b_6} \le M_t^{s_8} \le \dots \le M_t^{s_{12}} \le M_t^{b_1}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_T < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7

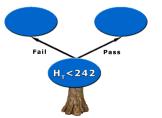




- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:

 $\begin{array}{l} p_{T}^{s_{1}} \leq p_{T}^{b_{24}} \leq \cdots \leq p_{T}^{b_{2}} \leq p_{T}^{s_{12}} \\ H_{T}^{b_{5}} \leq H_{T}^{b_{3}} \leq \cdots \leq H_{T}^{s_{67}} \leq H_{T}^{s_{43}} \\ M_{t}^{b_{6}} \leq M_{t}^{s_{8}} \leq \cdots \leq M_{t}^{s_{12}} \leq M_{t}^{b_{9}} \\ \end{array}$   $\begin{array}{l} \text{best split (arbitrary unit):} \\ p_{T} < 56 \text{ GeV, separation} = 3 \\ H_{T} < 242 \text{ GeV, separation} = 5 \\ \end{array}$ 

- $M_t < 105$  GeV, separation = 0.7
- split events in two branches: pass or fail *H*<sub>T</sub> < 242 GeV

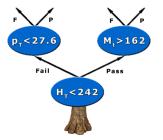




- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:  $p_{T1}^{s_1} \le p_{T}^{b_{24}} \le \dots \le p_{T}^{b_2} \le p_{T}^{s_{12}}$   $H_{T}^{b_5} \le H_{T}^{b_3} \le \dots \le H_{T}^{s_{67}} \le H_{T}^{s_{43}}$   $M_t^{b_6} \le M_t^{s_8} \le \dots \le M_t^{s_{12}} \le M_t^{b_9}$ best split (arbitrary unit):  $p_T < 56 \text{ GeV}$ , separation = 3  $H_T < 242 \text{ GeV}$ , separation = 5  $M_t < 105 \text{ GeV}$ , separation = 0.7
  - split events in two branches: pass or fail  $H_{\rm T} < 242$  GeV
- Repeat recursively on each node

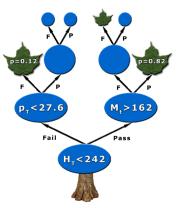






- Consider signal (s<sub>i</sub>) and background (b<sub>j</sub>) events described by 3 variables: p<sub>T</sub> of leading jet, top mass M<sub>t</sub> and scalar sum of p<sub>T</sub>'s of all objects in the event H<sub>T</sub>
  - sort all events by each variable:

    - $\blacksquare M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$
  - best split (arbitrary unit):
    - $p_T < 56$  GeV, separation = 3
    - $H_{\rm T} < 242$  GeV, separation = 5
    - $M_t < 105$  GeV, separation = 0.7
  - split events in two branches: pass or fail H<sub>T</sub> < 242 GeV</p>
- Repeat recursively on each node
- Splitting stops: e.g. events with  $H_{\rm T} < 242$  GeV and  $M_t > 162$  GeV are signal like (p = 0.82)





# **Decision 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

p-0.12 p-0.12 p-27.6 M\_2 162 Pass H\_1<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. <sup>1</sup>/<sub>2</sub>) in leaf
- E.g. events with H<sub>T</sub> < 242 GeV and M<sub>t</sub> > 162 GeV have a DT output of 0.82 or +1

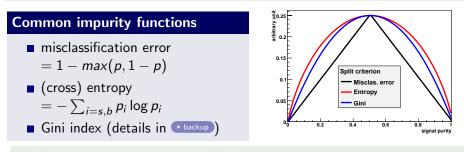
Yann Coadou (CPPM) — Machine learning

# Splitting a node



#### **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): ∆i(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  $\left(-\frac{s^2}{b}\right)$ 



#### 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:
  - $\blacksquare$  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):
  - $\blacksquare \ convert \ MeV \to GeV$
  - 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):
  - convert  $MeV \rightarrow GeV$
  - 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 selection III

#### Variable ranking (mean decrease impurity MDI)

- **•** Ranking of  $x_i$ : add up decrease of impurity each time  $x_i$  is used
- Largest decrease of impurity = best variable

#### Shortcoming: masking of variables

- $\mathbf{x}_i$  may be just a little worse than  $x_i$  but will never be picked
- $\blacksquare$  x<sub>i</sub> is ranked as irrelevant
- But remove  $x_i$  and  $x_i$  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

Yann Coadou (CPPM) — Machine learning

ESIPAP'22, Online, 26 Jan 2022 42/132

[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
- **B**rute 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.



- 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 ±1 if not using purity)

# Pruning a tree

#### 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 (details in • backup)

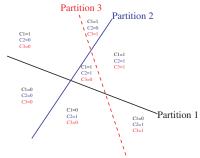
- Pre-pruning (early stopping condition like min leaf size, max depth)
- Expected error pruning (based on statistical error estimate)
- Cost-complexity pruning (penalise "complex" trees with many nodes/leaves)

# Tree (in)stability: distributed representation



#### One tree:

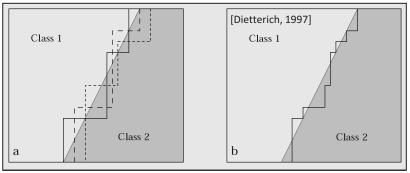
- one information about event (one leaf)
- cannot really generalise to variations not covered in training set (at most as many leaves as input size)
- Many trees:
  - distributed representation: number of intersections of leaves exponential in number of trees
  - $\blacksquare$  many leaves contain the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



# Tree (in)stability solution: averaging



#### Build several trees and average the output



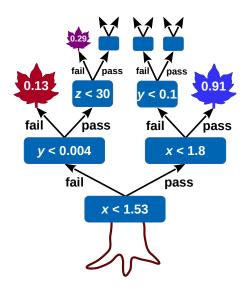
K-fold cross-validation (good for small samples)

- 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  $\mathcal{L}_k$
- **D**T output  $= \frac{1}{K} \sum_{k=1..K} T_k$

Bagging, boosting, random forests: ensemble learning

# **Boosted decision trees**





- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Kernel density estimation
- Decision trees
- Boosted decision trees
  - Introduction
  - AdaBoost
  - Clues to boosting performance
  - Gradient boosting
  - Performance examples
  - BDTs in real physics cases
- Neural networks
- Deep neural networks

# CPPM

### 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 $(T_1, T_2, T_3)$

# CPPM

#### 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 $(T_1, T_2, T_3)$

#### 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)

#### 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 $(T_1, T_2, T_3)$

#### 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  $T_k$ :
  - If 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 isMisclassified<sub>k</sub>(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  $\mathbb{T}_{k+1}$ :

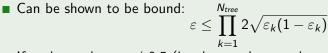
$$w_i^k \to w_i^{k+1} = w_i^k \times e^{\alpha_k}$$

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

# AdaBoost error rate



#### 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  $N_{\text{tree}}$ 

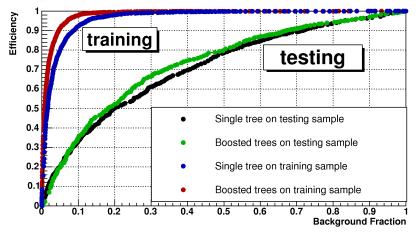
Corollary: training data is overfitted

#### **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

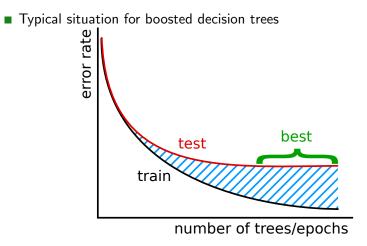


#### Efficiency vs. background fraction



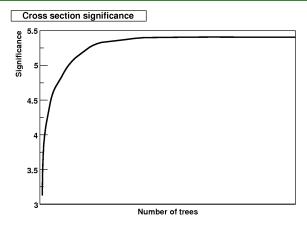
Clear overtraining, but still better performance after boosting

# Overtraining estimation: good or bad?



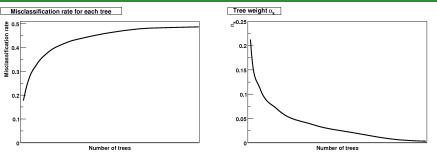
"bad" overtraining (overfitting) / "good" overtraining (still underfitting)

# 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

# Clues to boosting performance



- First tree is best, others are minor corrections
- Specialised trees do not perform well on most events ⇒ 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  $L_{MSF}(x, y) = \frac{1}{2} (y - F_k(x))^2$ 
  - 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:  $h_k(x_i) = y_i - F_k(x_i) = -\frac{\partial J}{\partial F_k(x_i)}$ 

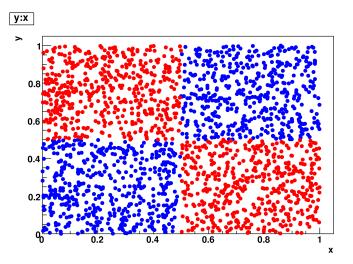
Generalised to any differentiable loss function



Friedman 2001

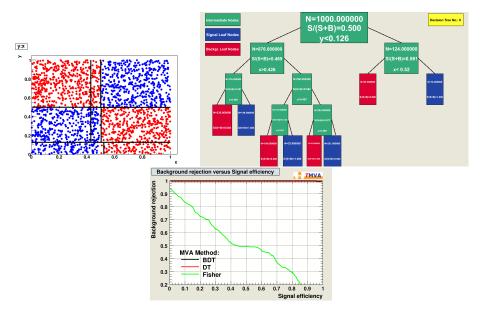
# Example: XOR problem





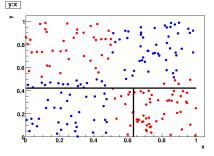
# Example: XOR problem





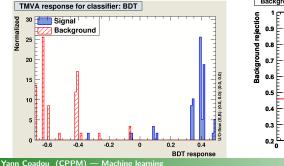
# Example: XOR with 100 events

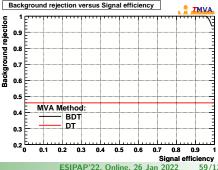




#### **Small statistics**

- Single tree not so good
- BDT very good: high performance discriminant from combination of weak classifiers

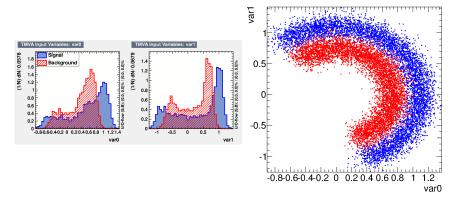




59/132

# **Circular correlation**

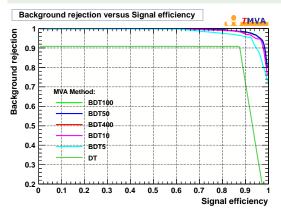




# **Circular correlation**

#### **Boosting longer**

- Compare performance of single DT and BDT with more and more trees (5 to 400)
- All other parameters unchanged

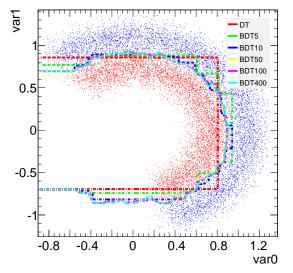


- Single (small) DT: not so good
- More trees ⇒ improve performance until saturation



# **Decision contours**

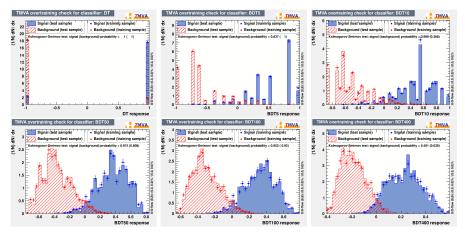




- Note: max tree depth = 3
- Single (small) DT: not so good. Note: a larger tree would solve this problem
- More trees ⇒ improve performance (less step-like, closer to optimal separation) until saturation
- Largest BDTs: wiggle a little around the contour ⇒ picked up features of training sample, that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

# Training/testing output

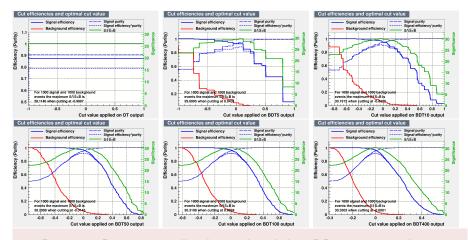




- Better shape with more trees: quasi-continuous
- Overtraining because of disagreement between training and testing? Let's see...

# Performance in optimal significance

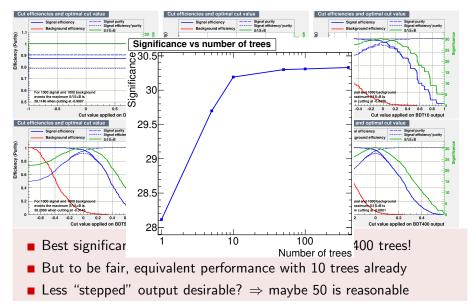


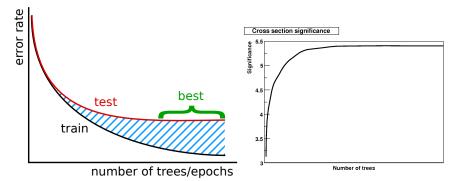


Best significance actually obtained with last BDT, 400 trees!
 But to be fair, equivalent performance with 10 trees already
 Less "stepped" output desirable? ⇒ maybe 50 is reasonable

# Performance in optimal significance









# Other averaging techniques

### 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<sub>k</sub> 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]

### **BDT in HEP**

### ATLAS b-tagging in Run 2

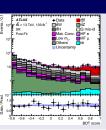
- Run 1 MV1c: NN trained from output of other taggers
- Run 2 MV2c20: BDT using feature variables of underlying algorithms and p<sub>T</sub>, η of jets
- Run 2: introduced IBL (new innermost pixel layer)
  - $\Rightarrow$  explains part of the performance gain, but not all

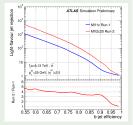
### ATLAS $t\bar{t}t\bar{t}$ production evidence

▶ 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

#### ESIPAP'22. Online. 26 Jan 2022 67/132



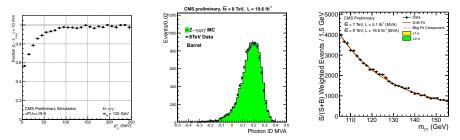




#### CMS-PAS-HIG-13-001

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

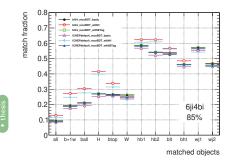


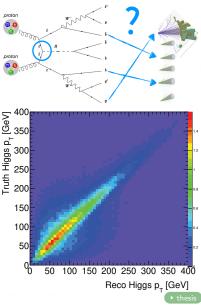
### **BDT in HEP: reducing combinatorics**



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

- Match jets and partons in high-multiplicity final state
- BDT trained on all combinations
- New inputs to classification BDT
- Access to Higgs p<sub>T</sub>, origin of b-jets
  Phys. Rev. D 97, 072016 (201)

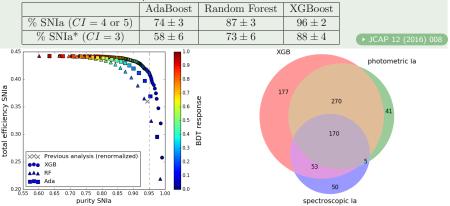




### BDT in cosmology: photometric classification



- Type la supernovae photometric classification
- Deriving redshifts from SN light curves on real data
- Tried random forest, AdaBoost and XGBoost



### Decision 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

### **Neural networks**





- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Kernel density estimation
- Decision trees
- Boosted decision trees
- Neural networks
  - History
  - Single neuron
  - Neural network training
- Deep neural networks

### **Neural networks**



### Neuron Human brain $\bullet$ 10<sup>11</sup> neurons 10<sup>14</sup> synapses Dendrites Learning: modifying synapses Ахоп Electrical Impulses Neurotransmitter Molecules Receptor Synapse



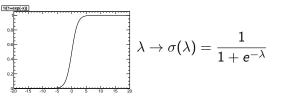
- 1943: W. McCulloch and W. Pitts explore capabilities of networks of simple neurons
- 1958: F. Rosenblatt introduces perceptron (single neuron with adjustable weights and threshold activation function)
- 1969: M. Minsky and S. Papert prove limitations of perceptron (linear separation only) and (wrongly) conjecture that multi-layered perceptrons have same limitations

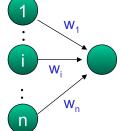
 $\Rightarrow$  ANN research almost abandoned in 1970s!!!

1986: Rumelhart, Hinton and Williams introduce "backward propagation of errors": solves (partially) multi-layered learning

### Single neuron

- Remember linear separation (Fisher discriminant): λ(x) = w ⋅ x = ∑<sup>n</sup><sub>i=1</sub> w<sub>i</sub>x<sub>i</sub> + w<sub>0</sub>
- Boundary at  $\lambda(x) = 0$
- Replace threshold boundary by sigmoid (or tanh):





- $\sigma$ : activation function (neuron activity)
- Neuron behaviour completely controlled by weights  $w = \{w_0, \ldots, w_n\}$
- Training: minimisation of error/loss function (quadratic deviations, entropy [maximum likelihood]), via gradient descent or stochastic approximation





#### Universal approximation theorem

Let  $\sigma(.)$  be a non-constant, bounded, and monotone-increasing continuous function. Let  $C(I_n)$  denote the space of continuous functions on the n-dimensional hypercube. Then, for any given function  $f \in C(I_n)$  and  $\varepsilon > 0$  there exists an integer M and sets of real constants  $w_j$ ,  $w_{ij}$  where i = 1, ..., n and j = 1, ..., M such that

$$y(x,w) = \sum_{j=1}^{M} w_j \sigma \left( \sum_{i=1}^{n} w_{ij} x_i + w_{0j} \right)$$

is an approximation of f(.), that is  $|y(x) - f(x)| < \varepsilon$ .

#### Interpretation

- You can approximate any continuous function to arbitrary precision with a linear combination of sigmoids
- Corollary 1: can approximate any continuous function with neurons!
- Corollary 2: a single hidden layer is enough
- Corollary 3: a linear output neuron is enough



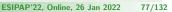
#### Interpretation

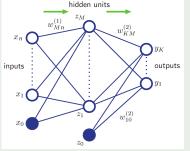
- You can approximate any continuous function to arbitrary precision with a linear combination of sigmoids
- Corollary 1: can approximate any continuous function with neurons!
- Corollary 2: a single hidden layer is enough
- Corollary 3: a linear output neuron is enough

#### Multilayer perceptron: feedforward network

- Neurons organised in layers
- Output of one layer becomes input to next layer

$$y_k(x,w) = \sum_{j=0}^{M} w_{kj}^{(2)} \underbrace{\sigma\left(\sum_{i=0}^{n} w_{ji}^{(1)} x_i\right)}_{z_i}$$









Training means minimising error function E(w)

•  $\frac{\partial E}{\partial w_j} = \sum_{n=1}^{N} - (t^{(n)} - y^{(n)}) x_j^{(n)}$  with target  $t^{(n)}$  (0 or 1), so  $t^{(n)} - y^{(n)}$  is the error on event n

- All events at once (batch learning):
  - weights updated all at once after processing the entire training sample
  - finds the actual steepest descent
  - takes more time
  - usually: mini-batches (send events by batches)
  - new training events: need to restart training from scratch
- or one-by-one (online learning):
  - incremental learning: new training events included as they come
  - speeds up learning
  - may avoid local minima with stochastic component in minimisation
  - careful: depends on the order of training events
- One epoch: going through the entire training data once

### Backpropagation

- Training means minimising error function *E(w)*
- For single neuron:  $\frac{dE}{dw_k} = (y t)x_k$

• One can show that for a network:

$$\frac{dE}{dw_{ji}} = \delta_j z_i, \text{ where }$$

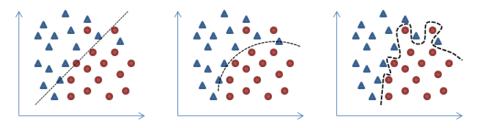
$$\delta_k = (y_k - t_k)$$
 for output neurons  
 $\delta_j \propto \sum_k w_{kj} \delta_k$  otherwise  
Hence errors are propagated backwards

$$z_i \bigcirc \delta_j \\ w_{ji} \bigcirc z_j \\ \delta_1 \\ \delta_1$$



### Neural network overtraining





Diverging weights can cause overfitting

#### Mitigate by:

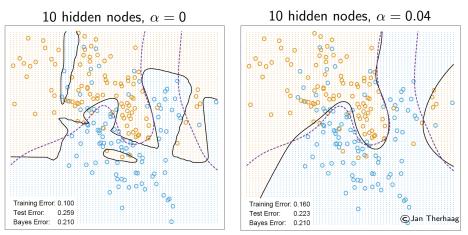
early stopping (after a fixed number of epochs)

- monitoring error on test sample
- regularisation, introducing a "weight decay" term to penalise large weights, preventing overfitting. For instance L2-regularisation:

$$\tilde{E}(w) = E(w) + \frac{\alpha}{2} \sum_{i} w_i^2$$

### Regularisation





Much less overfitting, better generalisation properties

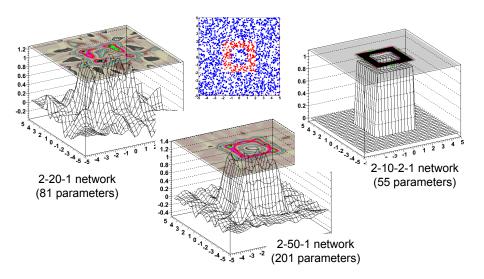


#### Preprocess data:

- if relevant, provide e.g. x/y instead of x and y
- subtract the mean because the sigmoid derivative becomes negligible very fast (so, input mean close to 0)
- normalise variances (close to 1)
- shuffle training sample (order matters in online training)
- Initial random weights should be small to avoid saturation
- Regularise weights to minimise overtraining
- Make sure the training sample covers the full parameter space
- No rule (not even guestimates) about the number of hidden nodes (unless using constructive algorithm, adding resources as needed)
- A single hidden layer is enough for all purposes, but multiple hidden layers may allow for a solution with fewer parameters

### Adding a hidden layer





### Deep neural networks



- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Kernel density estimation
- Decision trees
- Boosted decision trees
- Neural networks
- Deep neural networks
  - ANN interesting again!
  - New activation functions
  - Convolutional networks
  - Recurrent neural networks
  - Auto-encoders
  - Domain adaptation
  - Examples



### **Deep learning**



#### What is learning?

- Ability to learn underlying and previously unknown structure from examples
  - $\Rightarrow$  capture variations
- **Deep learning:** have several hidden layers (> 2) in a neural network

#### Motivation for deep learning

- Inspired by the brain (esp. visual cortex)
- Humans organise ideas hierarchically, through composition of simpler ideas
- Heavily unsupervised training, learning simpler tasks first, then combining into more abstract ones
- Learn first order features from raw inputs, then patterns in first order features, then etc.

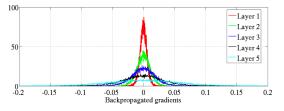
#### Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - "vanishing gradient": gradients getting very small further away from output ⇒ early layers do not learn much, can even penalise overall performance



#### Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - "vanishing gradient": gradients getting very small further away from output ⇒ early layers do not learn much, can even penalise overall performance





#### Deep networks were unattractive

- One layer theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
  - requires labelled data (usually scarce and expensive)
  - does not scale well, getting stuck in local minima
  - "vanishing gradient": gradients getting very small further away from output ⇒ early layers do not learn much, can even penalise overall performance

#### Breakthroughs around 2006 (Bengio, Hinton, LeCun)

- Train each layer independently (first practical approach)
- Can use unlabelled data (a lot of it), with unsupervised training
- New activation functions
- Possible thanks to algorithmic innovations, computing resources, data!

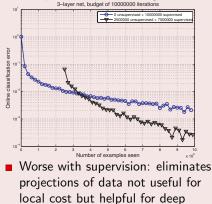


### Why does unsupervised training work?



#### Example

- Stacked denoising auto-encoders
- 10 million handwritten digits
- First 2.5 million used for unsupervised pre-training



### Optimisation hypothesis

- Training one layer at a time scales well
- Backpropagation from sensible features
- Better local minimum than random initialisation, local search around it

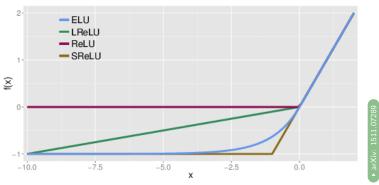
# Overfitting/regularisation hypothesis

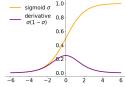
- More info in inputs than labels
- No need for final discriminant to discover features
- Fine-tuning only at category boundaries

model cost



- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

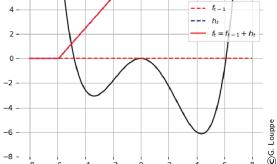


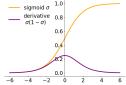


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.







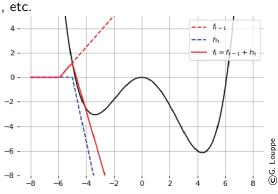


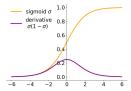


- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.





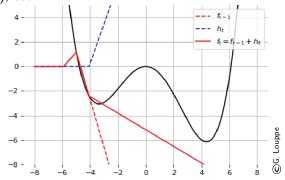


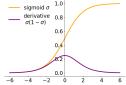




- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

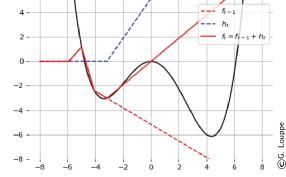


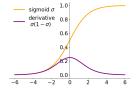






- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

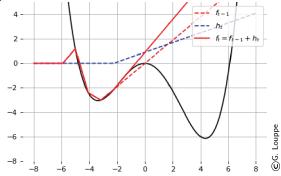


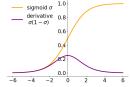






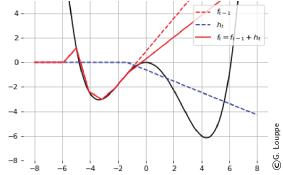
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

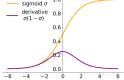


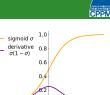




- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

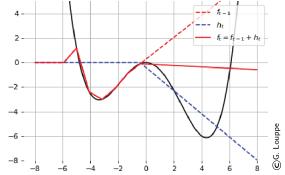






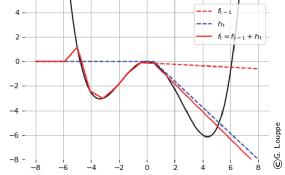
0.0

- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.





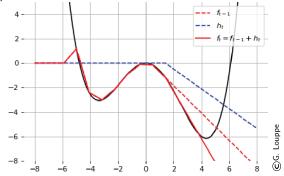
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

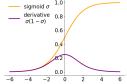


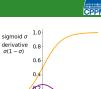




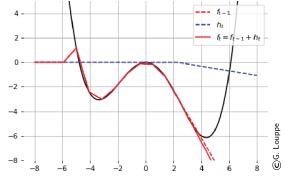
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

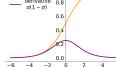






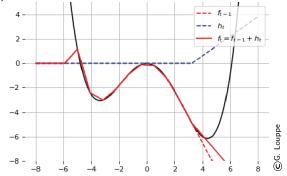
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

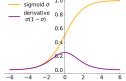






- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

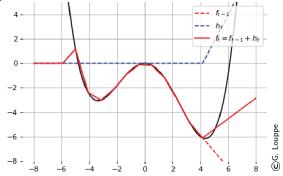


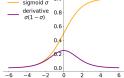


# Rectified linear unit (ReLU) activation



- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.

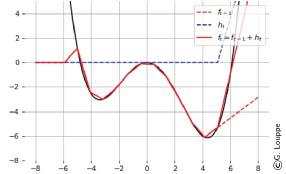


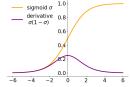


# Rectified linear unit (ReLU) activation



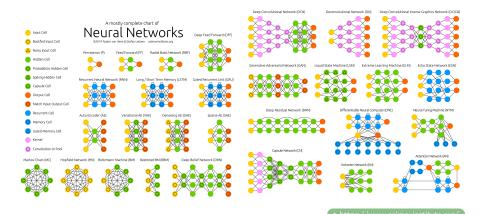
- One of reasons for vanishing gradient: sigmoid activation
  - tiny non-varying derivative away from zero
- Solution: non-saturating function
- Simplest case: rectified linear unit ReLU
- Other variants: leaky ReLU, shifted ReLU (SReLU), exponential linear unit (ELU), etc.





### Neural network zoo



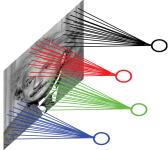


Many possible network structures

Moving away from feature engineering to model design



Images are stationary: can learn feature in one part and apply it in another





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

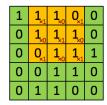




4		



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

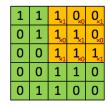


Image





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

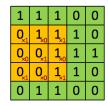


Image

4	3	4



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

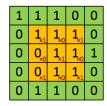


Image





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image



Image





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

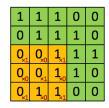
1	1	1	0	0
0	1	<b>1</b> <sub>×1</sub>	<b>1</b> _×0	<b>0</b> ,×1
0	0	<b>1</b> <sub>×0</sub>	<b>1</b> <sub>×1</sub>	1 <sub>×0</sub>
0	0	<b>1</b> _×1	1_×0	<b>0</b> ,×1
0	1	1	0	0

Image





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

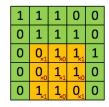




4	3	4
2	4	3
2		



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image



Image





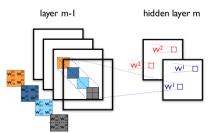
- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0	0	<b>1</b> <sub>×1</sub>	1_×0	1,
0	0	<b>1</b> <sub>×0</sub>	<b>1</b> _×1	<b>0</b> _×0
0	1	1,×1	<b>0</b> _×0	<b>0</b> _×1

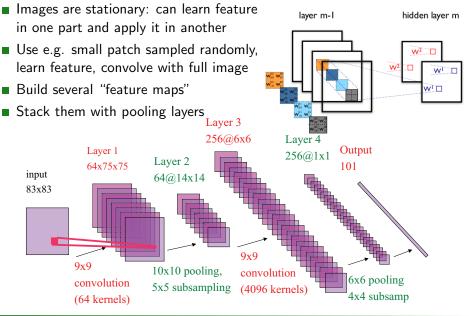
Image



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image
- Build several "feature maps"

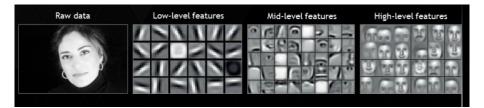










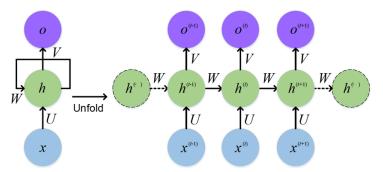


### Many problems require processing a sequence

- sequence classification
  - text analysis ("sentiment analysis")
  - DNA sequencing
  - action selection
- sequence synthesis
  - text synthesis
  - music/video
- sequence translation
  - speech recognition
  - translation
- Usually variable length sequences (number of words/ notes/ frames/ etc.)
- Use a recurrent model, maintaining a recurrent state updated after each step

### **Recurrent neural networks**

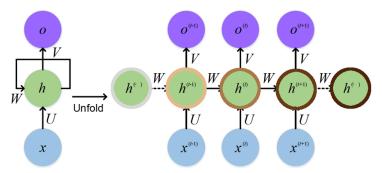




- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions

### **Recurrent neural networks**

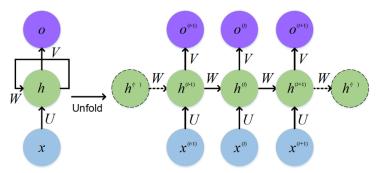




- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions
- Issue: early frames diluted over sequence ⇒ memory loss

### **Recurrent neural networks**

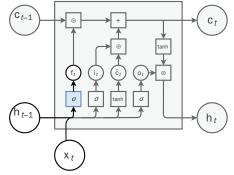




- Keeps information from earlier frames while processing (variable-size) sequence
- Could also be bi-directional, consuming sequence in both directions
- Issue: early frames diluted over sequence  $\Rightarrow$  memory loss
- Introducing long short-term memory (LSTM) networks
  - using forget gate to regulate information flow
  - also possible with gated recurrent units (GRU)

### Long short-term memory (LSTM)

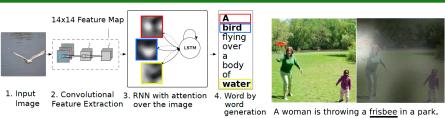




- Recurrent state split in two parts
  - cell state  $c_t$
  - output state  $h_t$
- Forget gate  $f_t$  to erase cell state info
- Input gate it to update cell state info
- Output gate o<sub>t</sub> to select output state

### Recurrent neural networks examples

### Labelling images

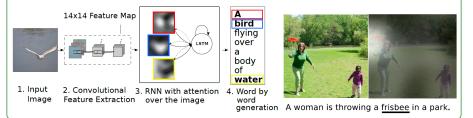




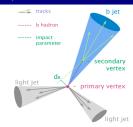
### Recurrent neural networks examples

### Labelling images





### b-jet tagging in ATLAS experiment



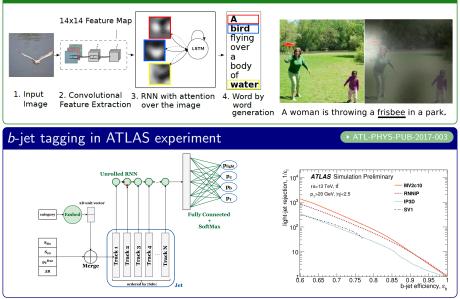
### • ATL-PHYS-PUB-2017-003

### Recurrent neural networks examples

### Labelling images



▶ arXiv:1502.03044



### **Auto-encoders**

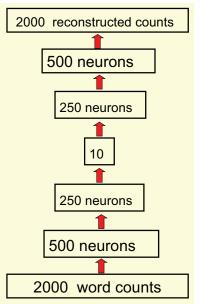


### Approximate the identity function

- Build a network whose output is similar to its input
- Sounds trivial? Except if imposing constraints on network (e.g., # of neurons, locally connected network) to discover interesting structures
- Can be viewed as lossy compression of input

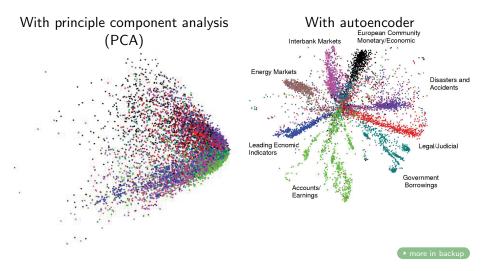
### Finding similar books

- Get count of 2000 most common words per book
- "Compress" to 10 numbers



### **Auto-encoders**

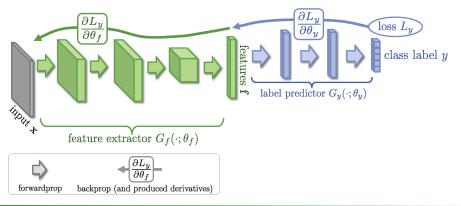




### Domain adaptation and adversarial training

### Typical training

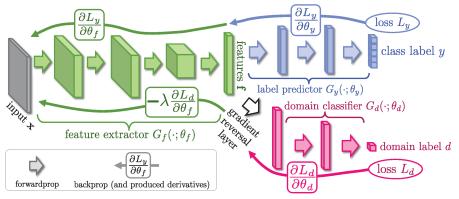
- signal and background from simulation
- results compared to real data to make measurement
- Requires good data-simulation agreement



### Domain adaptation and adversarial training

### Typical training

- signal and background from simulation
- results compared to real data to make measurement
- Requires good data-simulation agreement
- Possibility to use adversarial training and domain adaptation to account for discrepancies/systematic uncertainties





### Deep networks and art

- Learning a style arXiv:1508.06576 [cs.CV] Neural-style
   Image: ArXiv:1508.06576 [cs.CV] Neural-style

Computer dreams • Google original

deepdream





Face Style • http://facestyle.org

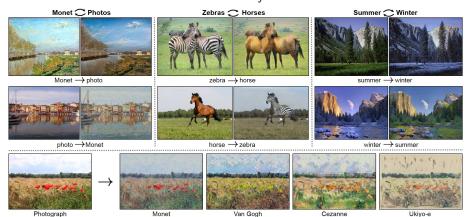
http://dcgi.fel.cvut.cz/home/sykorad/facestyle.html



### CycleGAN



### ■ GAN = generative adversarial network Summary



### CycleGAN



### $\blacksquare \ \mathsf{GAN} = \mathsf{generative} \ \mathsf{adversarial} \ \mathsf{network}$



### CycleGAN





### ILSVRC



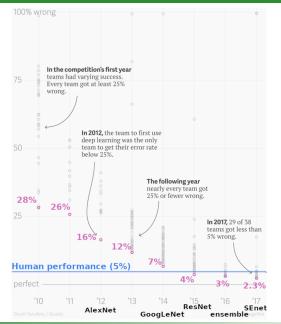
### ImageNet Large Scale Visual Recognition Challenge

- ImageNet: database with 14 million images and 20k categories
- Used 1000 categories and about 1.3 million manually annotated images



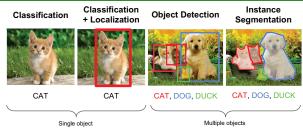
### **ILSVRC** classification performance





### **Increasing refinement**





- More and more refinement (segmentation)
- More objects, in real time on video1/video2/video3





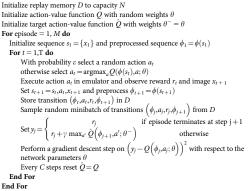
→ Nature 518, 529 (2015)

- Learning to play 49 different Atari 2600 games
- No knowledge of the goals/rules, just 84x84 pixel frames
- 60 frames per second, 50 million frames (38 days of game experience)
- Deep convolutional network with reinforcement: DQN (deep Q-network)
  - action-value function  $Q^*(s,a) = \max_{\pi} \mathbb{E} \left[ r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi \right]$
  - maximum sum of rewards  $r_t$  discounted by  $\gamma$  at each timestep t, achievable by a behaviour policy  $\pi = P(a|s)$ , after making observation s and taking action a
- Tricks for scalability and performance:
  - experience replay (use past frames)
  - separate network to generate learning targets (iterative update of Q)
- Outperforms all previous algorithms, and professional human player on most games

# Google DeepMind: training&performance



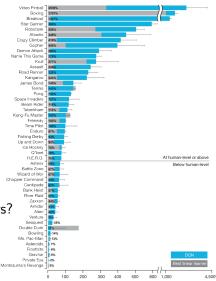
#### Algorithm 1: deep Q-learning with experience replay.



#### • What about Breakout or Space invaders?



	0P30 0000
	<b>КХХХХ</b> Х
*** **	
	*****
	*****
22 22	888888
	****
A A A	A A'A
	4



### Google DeepMind: mastering Go



- Board games: can be solved with search tree of b<sup>d</sup> possible sequences of moves
   b = breadth [number of legal moves]
  - *d* = depth [length of game]

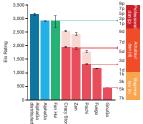


▶ Nature 529, 484 (2016)

- Chess:  $b \approx 35$ ,  $d \approx 80 \rightarrow$  go:  $b \approx 250$ ,  $d \approx 150$
- Reduction:
  - of depth by position evaluation (replace subtree by approximation that predicts outcome)
  - of breadth by sampling actions from probability distribution (policy p(a|s)) over possible moves a in position s
- $19 \times 19$  image, represented by CNN
- Supervised learning policy network from expert human moves, reinforcement learning policy network on self-play (adjusts policy towards winning the game), value network that predicts winner of games in self-play.

# Google DeepMind: AlphaGo

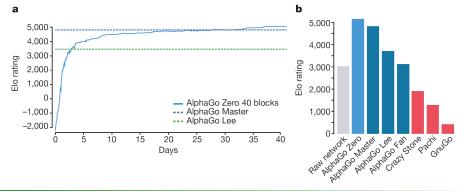
- AlphaGo: 40 search threads, simulations on 48 CPUs, policy and value networks on 8 GPUs. Distributed AlphaGo: 1020 CPUs, 176 GPUs
- AlphaGo won 494/495 games against other programs (and still 77% against Crazy Stone with four handicap stones)
- Fan Hui: 2013/14/15 European champion
- Distributed AlphaGo won 5–0
- AlphaGo evaluated thousands of times fewer positions than Deep Blue (first chess computer to bit human world champion) ⇒ better position selection (policy network) and better evaluation (value network)



- Then played Lee Sedol (top Go play in the world over last decade) in March 2016 ⇒ won 4–1. AlphaGo given honorary professional ninth dan, considered to have "reach a level 'close to the territory of divinity' "
- Ke Jie (Chinese world #1): "Bring it on!". May 2017: 3–0 win for AlphaGo. New comment: "I feel like his game is more and more like the 'Go god'. Really, it is brilliant"

### DeepMind AlphaGo Zero

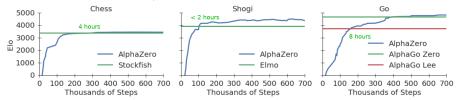
- Learn from scratch, just from the rules and random moves
- Reinforcement learning from self-play, no human data/guidance
- Combined policy and value networks
- 4.9 million self-play games
- Beats AlphaGo Lee (several months of training) after just 36 hours
- Single machine with four TPU



▶ Nature 550, 354 (2017)



- Same philosophy as AlphaGo Zero, applied to chess, shogi and go
- Changes:
  - not just win/loss, but also draw or other outcomes
  - no additional training data from game symmetries
  - using always the latest network to generate self-play games rather than best one
  - tree search: 80k/70M for chess AlphaZero/Stockfish, 40k/35M for shogi AlphaZero/Elmo

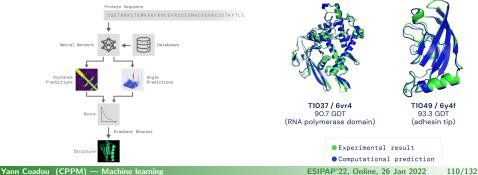


## DeepMind AlphaFold

- Trying to tackle scientific problem
- Goal: predict 3D structure of protein based solely on genetic sequence
- Using DNN to predict
  - distances between pairs of amino acids
  - angles between chemical bonds
- Search DB to find matching existing substructures
- Also train a generative NN to invent new fragments
- Achieved best prediction ever



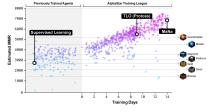
→ Blog Dec 2018 → AlphaFold2 Jul 2021

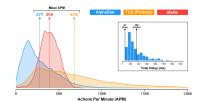


#### Yann Coadou (CPPM) — Machine learning

#### DeepMind AlphaStar

- Mastering real-time strategy game StarCraft II
- Challenges in game theory (no single best strategy), imperfect information (hidden parts of game), long term planning, real time (continuous flow of actions), large action space (many units/buildings)
- Using DNN trained
  - directly on raw data games
  - supervised learning on human games
  - reinforcement learning (continuous league)
- DNN output: list of actions
- Trained for 14 days; each agent: up to 200 years of real-time play
- Runs on single desktop GPU
- Defeated 5–0 one of best pro-players







▶ Blog Jan 2019 ▶ Nature 575, 350–354 (2019)



#### Deep networks at work



#### Playing poker

- Libratus (AI developed by Carnegie Mellon University) defeated four of the world's best professional poker players (Jan 2017 • arXiv:1705.02955))
- After 120,000 hands of Heads-up, No-Limit Texas Hold'em, led the pros by a collective \$1,766,250 in chips
- Learned to bluff, and win with incomplete information and opponents' misinformation
- Lip reading arXiv:1611.05358 [cs.CV]
  - human professional: deciphers less than 25% of spoken words
  - CNN+LSTM trained on television news programs: 50%

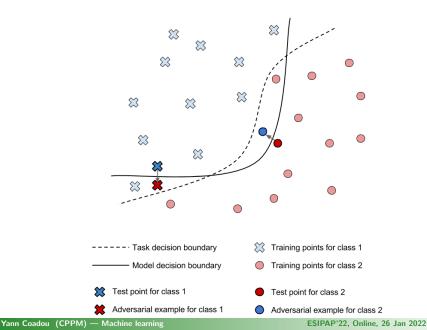
Limitation: adversarial attacks • arXiv:1312.6199 [cs.CV]



- left: correctly classified image
- middle: difference between left image and adversarial image (x10)
- right: adversarial image, classified as ostrich

#### Adversarial attack: what is happening?









original semantic segmentation framework





original semantic segmentation framework



adversarial attack





original semantic segmentation framework



adversarial attack



compromised semantic segmentation framework



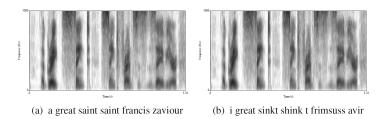
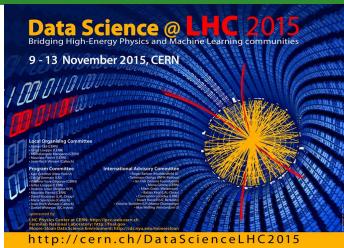


Figure 7: The model models' output for each of the spectrograms is located at the bottom of each spectrogram. The target transcription is: A Great Saint Saint Francis Xavier.





http://opendata.cern.ch

Yann Coadou (CPPM) — Machine learning



Data Science @ LHC 2015 Bridging High-Energy Physics and Machine Learning communities

Exploring the potential for Machine Learning on ATLAS

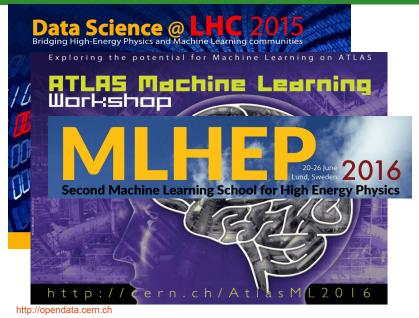
## ATLAS Machine Learning Workshop

29th-31st March 2016, CERN

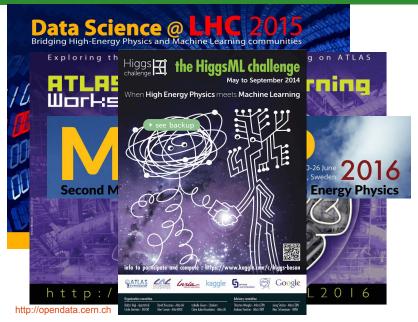
Organising Committee: Matthew Beckingham (Warwick) Michael Kagan (SLAC) David Rousseau (LAL-Orsav)

http://cern.ch/AtlasML2016

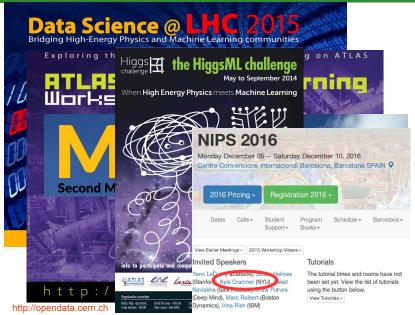








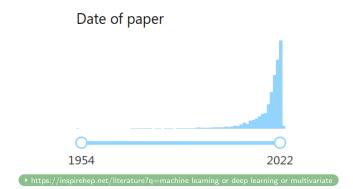












#### Up-to-date comprehensive review of papers

https://github.com/iml-wg/HEPML-LivingReview



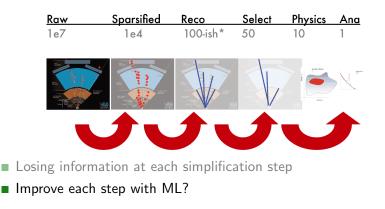
Reduce data dimensionality to allow analysis

Raw	Sparsified	Reco	Select	Physics	Ana
1e7	1e4	100-ish *	50	10	1
				and a second	A more
- <u> </u>   / /					X
and a second					
	- <b>19</b>	- TEY ::		LATAS .	

#### Losing information at each simplification step

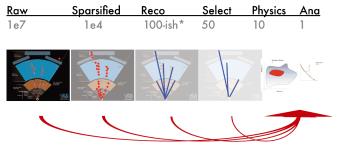


Reduce data dimensionality to allow analysis



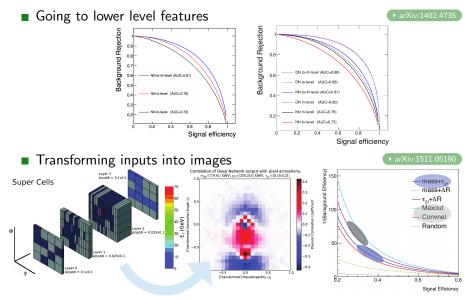


Reduce data dimensionality to allow analysis



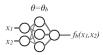
- Losing information at each simplification step
- Improve each step with ML?
- Skip one or more steps with ML?

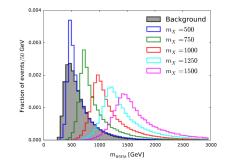






 ■ Looking for new physics scenario with unknown mass
 ⇒ one NN for each mass point θ=θ<sub>n</sub>







Background  $m_{Y} = 500$ 

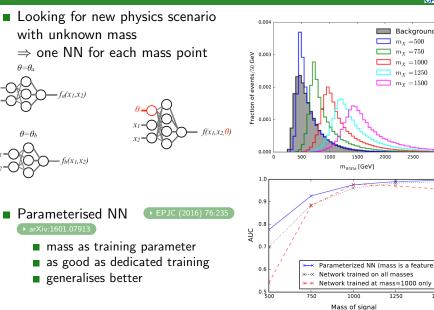
 $m_{Y} = 750$ 

 $m_{X} = 1000$ 

 $m_X = 1250$  $m_{X} = 1500$ 

2500

2000



Yann Coadou (CPPM) — Machine learning

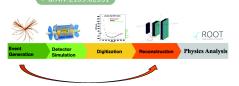
1250

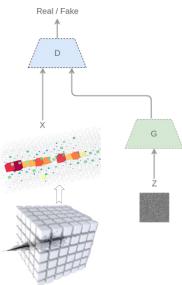
1500



Fast simulation with generative models

- Heavy CPU cost of simulation (> 50% of grid resources)
  - MC stats becoming limiting factor in analyses
- Replace "full simulation" with approximation
  - already routinely done, using parameterisation of showers or library of pre-simulated objects
  - use GAN to simulate medium-range hadrons in ATLAS





## Machine learning and systematics

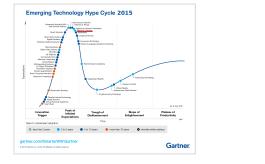


- No particular rule
- ML algorithm output can be considered as any other cut variable (just more powerful). Evaluate systematics by:
  - varying cut value
  - retraining
  - calibrating, etc.
- Most common (and appropriate): propagate other uncertainties (detector, theory, etc.) up to ML algorithm ouput and check how much the analysis is affected
- More and more common: profiling Watch out:
  - ML algorithm output powerful
  - signal region (high ML algorithm 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-ML selection cuts)



- Very active field of research in machine learning and artificial intelligence
  - not just at universities (Google, Facebook, Microsoft, NVIDIA, etc...)
- Training with curriculum:
  - what humans do over 20 years, or even a lifetime
  - learn different concepts at different times
  - solve easier or smoothed version first, and gradually consider less smoothing
  - exploit previously learned concepts to ease learning of new abstractions
- Combination of deep learning and reinforcement learning
- Domain adaptation and adversarial training
  - e.g. train in parallel network that produces difficult examples
  - learn discrimination (s vs. b) and difference between training and application samples (e.g. Monte Carlo simulation and real data)
- Getting popular: graph networks

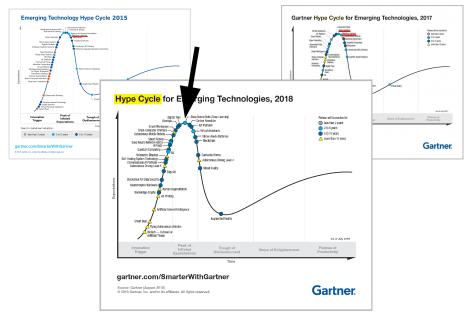




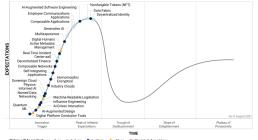


Hype cycle









Plateau will be reached: 🔿 < 2 vrs. 💿 2-5 vrs. 🌰 5-10 vrs. 🔺 >10 vrs. 🛠 Obsolete before plateau

Source: Gartner (August 2021)

No deep learning/ML anymore (since 2019)

Instead, all sorts of "AI" and ML-driven systems

> data fabric

- Al-augmented software engineering
- generative AI
- physics-informed AI
- Al-driven innovation
- Al-augmented design
- quantum ML

## Hype cycle 2021



#### Hype Cycle for Artificial Intelligence, 2021



Sources Gather 0.2027 Gather, Inc. and/or its efficients. All rights reserved. Gather and Hype Cycle are registered trademarks of Gather, Inc. and its efficients in the U.S. 1482644 Gartner.

- No deep learning/ML anymore (since 2019)
- Instead, all sorts of "AI" and ML-driven systems

data fabric

- Al-augmented software engineering
- generative AI
- physics-informed AI
- Al-driven innovation
- Al-augmented design
- quantum ML

## AI and society

- ML achieves super-human performance for well-designed problems, or games with score ⇒ where one can define a proper loss function or reward
- Scale to "real" problems?
  - explainability
  - causality
  - "moral" stand
  - culture, art
- Many advances in medical imaging, modelling of various phenomena, supernova analysis or LHC physics, but issues with:
  - out-of-distribution generalisation
  - scalability of computing resources, carbon footprint
  - reliability
  - decision bias (gender, race, etc.)
- Workshops/socials on Fairness & ethics, AI for Good, Tackling Climate Change with ML, AI for Humanitarian Assistance and Disaster Response, Safety and Robustness in Decision-making, ...
- Importance of personal decisions



#### From NeurIPS2019: Hidden information



#### Predicting the Politics of an Image Using Webly Supervised Data

Original	Reconst.	Far left	Far right	Original	Reconst.	Far left	Far right	Original	Reconst.	Far left	Far right
an	35	38	ap.	60	1	-	90		20		70
(4)	0	S	0	-	10	0	1		5,	0	1=
Ind.	N.	050	T	10.90	6.11	50	P	to a	0.0	100	2º
1 TO	104	00.0	1	(C ))	and a	No.	No.	E	T	10.0	T

#### Computational Mirrors: Blind Inverse Light Transport by Deep Matrix Factorization

#### (edited video)



#### Conclusion



- Many techniques and tools exist to achieve optimal discrimination
- (Un)fortunately, no one method can be shown to outperform the others in all cases
- One should try several and pick the best one for any given problem
- Latest machine learning algorithms (e.g. deep networks) require enormous hyperparameter space optimisation...
- Machine learning and multivariate techniques are at work in your everyday life without your knowning and can easily outsmart you for many tasks
- Try this to convince yourself http://www.phi-t.de/mousegame/mousegame\_en.html

#### Upcoming reference book (out next month)

Artificial Intelligence for High Energy Physics

https://doi.org/10.1142/12200

27/132

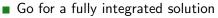
Forthcoming ARTIFICIAL INTELLIGENCE HIGH ENERGY

PHYSICS

Raolo Calafiura - David Bausseau - Kazuhiro Tera

### Machine learning software

CPPM



- 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**. See tutorial next week
  - 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]
- Specifically for Deep Learning:
  - TensorFlow (+Keras) [Google]
  - PyTorch [Facebook]
  - Torch (in lua)

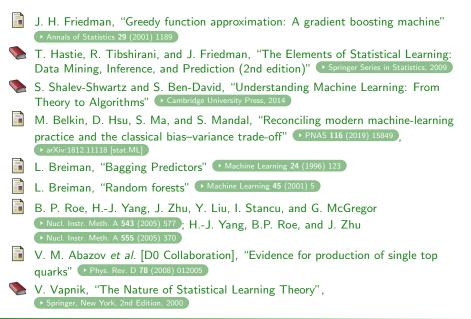
https://github.com/dmlc/xgboost



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

### References II: decision trees and more



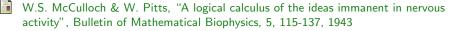






📎 C.M. Bishop, Pattern Recognition and Machine Learning, Springer, New York, 2007





F. Rosenblatt, "The Perceptron: A Probabilistic Model for Information Storage & Organization in the Brain", Psychological Review, 65, pp. 386-408, 1958



K.Hornik et al., "Multilayer Feedforward Networks are Universal Approximators", Neural Networks, Vol. 2, pp 359-366, 1989



Y. LeCun, L. Bottou, G. Orr and K. Muller, "Efficient BackProp", in Neural Networks: Tricks of the trade, Orr, G. and Muller K. (Eds), Springer, 1998

### References IV: deep neural networks



- Q.V. Le et al., "Building High-level Features Using Large Scale Unsupervised Learning", in *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012 • http://research.google.com/pubs/pub38115.html
- Y. Bengio, P. Lamblin, D. Popovici and H. Larochelle, "Greedy Layer-Wise Training of Deep Networks", in *Advances in Neural Information Processing Systems 19* (NIPS'06), pages 153–160, MIT Press 2007
- M.A. Ranzato, C. Poultney, S. Chopra and Y. LeCun, in J. Platt et al., "Efficient Learning of Sparse Representations with an Energy-Based Model", in *Advances in Neural Information Processing Systems 19* (NIPS'06), pages 1137–1144, MIT Press, 2007
- Y. Bengio, "Learning deep architectures for AI", Foundations and Trends in Machine Learning, Vol. 2, No. 1 (2009) 1–127. Also book at
   http://www.iro.umontreal.ca/lisa/publications2/index.php/publications/show/239



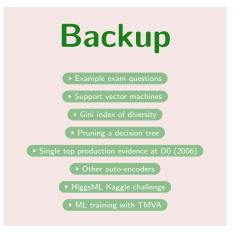
I. Goodfellow, Y. Bengio and A. Courville, *Deep Learning*, MIT Press, 2016 http://www.deeplearningbook.org



G. Carleo et al., "Machine learning and the physical sciences"

▶ Rev. Mod. Phys. 91, 045002 (2019) ▶ arXiv:1903.10563



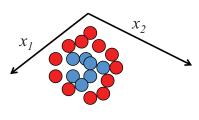


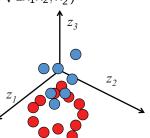


- Explain in your own words the three main approaches to train machine learning algorithms
- Comment on differences between boosted decision trees and neural networks: input variables, resilience to overtraining, hyperparameter tuning
- Explain why neural networks used to be shallow (mathematical motivation, training)
- Describe briefly the main features and use cases of CNN, RNN and auto-encoders
- Explain briefly the behaviour on each of the four plots on page 18, showing the training and testing curves as a function of the number of trees or epochs (difference between train/test, overfitting/underfitting, optimal classifier, etc.)

### Support vector machines

- 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:  $h: (x_1, x_2) \rightarrow (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$



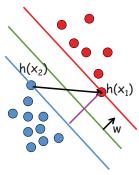


[Vapnik 2000]



• 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:  $R(w) \propto ||w||^2$

### Support vector machines: constraints



- 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 ·

 $w \cdot h(x) + b > 1$  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:

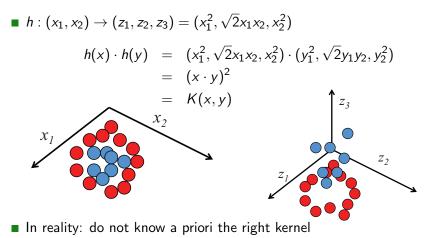
- need to find h mappings (potentially of infinite dimension)
- need to compute scalar products  $h(x_i) \cdot h(x_j)$
- Fortunately  $h(x_i) \cdot h(x_j)$  are equivalent to some kernel function  $K(x_i, x_j)$  that does the mapping and the scalar product:

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



### Support vector machines: example





 $\blacksquare$   $\Rightarrow$  have to test different standard kernels and use the best one

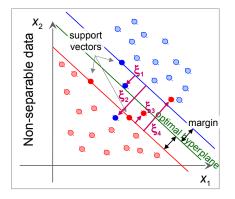
### Support vector machines: non-separable data



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

• 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<sub>i</sub>.
- Probability that it is actually in class j is p<sub>j</sub>
- $\blacksquare \Rightarrow \mathsf{Gini} = \mathsf{probability} \text{ of misclassification}$

### 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

## Pruning a tree I

### **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  $N_T$  nodes) of fully grown tree  $T_{max}$ :

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,  $T_{max}$  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 III: cost-complexity pruning



- For node t and subtree  $T_t$ :
  - if t non-terminal,  $R(t) > R(T_t)$  by construction
  - $\blacksquare R_{\alpha}(\lbrace t \rbrace) = R_{\alpha}(t) = R(t) + \alpha \ (N_T = 1)$
  - 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:  $\rho_t = \frac{R(t) R(T_t)}{N_{\tau} 1}$
  - **I** 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

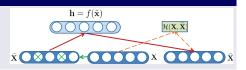
Note: best pruned tree may not be optimal in a forest

### Sparse auto-encoder

- Sparsity: try to have low activation of neurons (like in the brain)
- Compute average activation of each hidden unit over training set
- Add constraint to cost function to make average lower than some value close to 0

### Denoising auto-encoder

- Stochastically corrupt inputs
- Train to reconstruct uncorrupted input



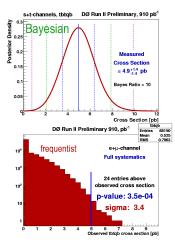
### Locally connected auto-encoder

- Allow hidden units to connect only to small subset of input units
- Useful with increasing number of input features (e.g., bigger image)
- Inspired by biology: visual system has localised receptive fields

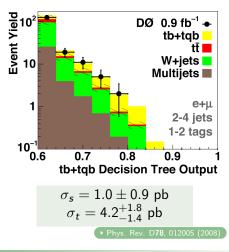
# Single top production evidence at D0 (2006)



- Three multivariate techniques: BDT, Matrix Elements, BNN
- Most sensitive: BDT



 $\sigma_{s+t} = 4.9 \pm 1.4 \text{ pb}$ p-value = 0.035% (3.4 $\sigma$ ) SM compatibility: 11% (1.3 $\sigma$ )



### Decision trees — 49 input variables



#### **Object Kinematics**

 $\begin{array}{l} p_{T}(jet1) \\ p_{T}(jet2) \\ p_{T}(jet3) \\ p_{T}(jet4) \\ p_{T}(otbest1) \\ p_{T}(notbest1) \\ p_{T}(notbest2) \\ p_{T}(tag1) \\ p_{T}(untag1) \\ p_{T}(untag2) \end{array}$ 

#### Angular Correlations

 $\Delta R(\text{iet1.iet2})$ cos(best1,lepton)besttop cos(best1,notbest1)besttop cos(tag1,alljets)<sub>alljets</sub> cos(tag1,lepton)btaggedtop cos(jet1,alljets)alljets cos(jet1,lepton)btaggedtop cos(jet2,alljets)<sub>alljets</sub> cos(jet2,lepton)btaggedtop  $\cos(\text{lepton}, Q(\text{lepton}) \times z)_{\text{besttop}}$ cos(lepton<sub>hesttop</sub>,besttop<sub>CMframe</sub>) cos(lepton<sub>btaggedtop</sub>, btaggedtop<sub>CMframe</sub>) cos(notbest,alljets)alliets cos(notbest,lepton) cos(untag1,alljets)<sub>alljets</sub> cos(untag1,lepton)

### **Event Kinematics**

Aplanarity(alljets, W) M(W.best1) ("best" top mass) M(W,tag1) ("b-tagged" top mass)  $H_{T}(alliets)$  $H_{T}(\text{alljets}-\text{best1})$  $H_{\rm T}(\text{alljets}-\text{tag1})$  $H_{T}(alljets, W)$  $H_{T}$ (jet1, jet2)  $H_{T}(jet1, jet2, W)$ M(alliets) M(alljets-best1) M(alliets-tag1) M(iet1.iet2) M(jet1, jet2, W) $M_{T}$ (jet1, jet2)  $M_{\rm T}(W)$ Missing  $E_{T}$ p<sub>T</sub>(alljets-best1) pT(alljets-tag1)  $p_{T}(jet1, jet2)$  $Q(lepton) \times \eta(untag1)$  $\sqrt{\hat{s}}$ Sphericity(alliets, W)

- Adding variables did not degrade performance
- Tested shorter lists, lost some sensitivity
- Same list used for all channels

### Decision trees — 49 input variables



#### **Object Kinematics**

 $\begin{array}{l} p_{T}(jet1) \\ p_{T}(jet2) \\ p_{T}(jet3) \\ p_{T}(jet4) \\ p_{T}(otbest1) \\ p_{T}(notbest1) \\ p_{T}(notbest2) \\ p_{T}(tag1) \\ p_{T}(untag1) \\ p_{T}(untag2) \end{array}$ 

#### Angular Correlations

 $\Delta R(\text{iet1.iet2})$ cos(best1,lepton)<sub>bestton</sub> cos(best1,notbest1)besttop cos(tag1,alljets)<sub>alljets</sub> cos(tag1,lepton)btaggedtop cos(jet1,alljets)alljets cos(jet1,lepton)btaggedtop cos(jet2,alljets)<sub>alljets</sub> cos(jet2,lepton)btaggedtop  $\cos(\text{lepton}, Q(\text{lepton}) \times z)_{\text{besttop}}$  $cos(lepton_{besttop}, besttop_{CMframe})$ cos(lepton<sub>btaggedtop</sub>, btaggedtop<sub>CMframe</sub>) cos(notbest,alljets)alliets cos(notbest,lepton) cos(untag1,alljets)<sub>alljets</sub> cos(untag1,lepton)

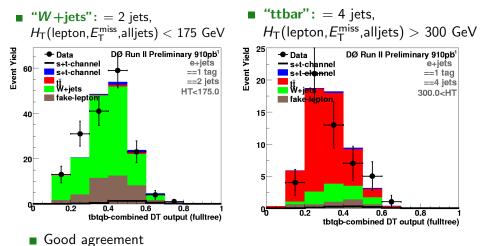
### **Event Kinematics**

Aplanarity(alljets, W) M(W.best1) ("best" top mass) M(W,tag1) ("b-tagged" top mass)  $H_{T}(alliets)$  $H_{T}(\text{alljets}-\text{best1})$  $H_{\rm T}(\text{alljets}-\text{tag1})$  $H_{T}(alljets, W)$  $H_{T}$ (jet1, jet2)  $H_{T}(jet1, jet2, W)$ M(alliets) M(alljets-best1) M(alliets-tag1) M(iet1.iet2) M(jet1, jet2, W) $M_{T}$ (jet1, jet2)  $M_{\rm T}(W)$ Missing  $E_{T}$ p<sub>T</sub>(alljets-best1) pT(alljets-tag1)  $p_{T}(jet1, jet2)$  $Q(lepton) \times \eta(untag1)$  $\sqrt{\hat{s}}$ Sphericity(alljets, W)

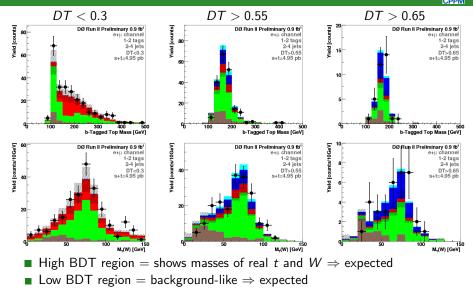
- Adding variables did not degrade performance
- Tested shorter lists, lost some sensitivity
- Same list used for all channels
- Best theoretical variable: *H*<sub>T</sub>(alljets, *W*). But detector not perfect ⇒ capture the essence from several variations usually helps "dumb" MVA



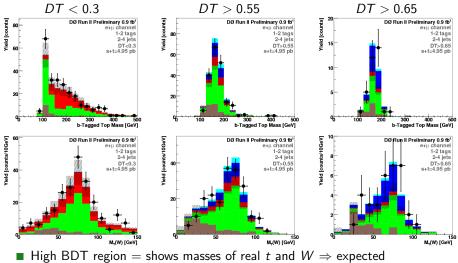
Validate method on data in no-signal region



### Boosted decision tree event characteristics



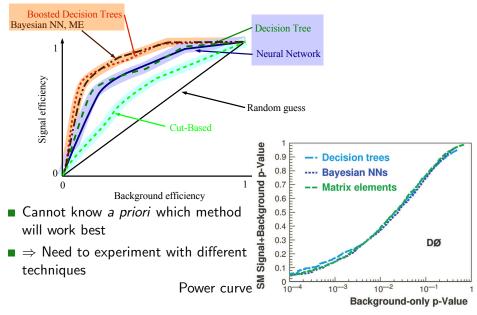
### Boosted decision tree event characteristics



- Low BDT region = background-like  $\Rightarrow$  expected
- Above does NOT tell analysis is ok, but not seeing this could be a sign of a problem

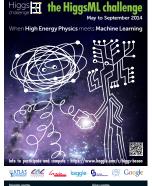
## Comparison for D0 single top evidence



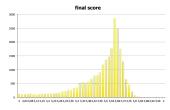


### HiggsML Kaggle challenge





Organization committee			Advisory committee		
		habelle Gapen - Chalkern Cleine Adam Baarderice - Aller CAE			



### HiggsML challenge

- Put ATLAS Monte Carlo samples on the web  $(H \rightarrow \tau \tau \text{ analysis})$
- Compete for best signal-bkg separation
- 1785 teams (most popular challenge ever)
- 35772 uploaded solutions

-	Se	e	<sup>ggle</sup> web site	and	mc	ore information
	∆rank '	Team Name amodel	uploaded * in the money	Score 🧐	Entries	Last Submission UTC (Best - Last Submission)
1	†1	Gábor Melis ‡ *	7000\$	3.80581	110	Sun, 14 Sep 2014 09:10:04 (-0h)
2	11	Tim Salimans ‡		3.78913	57	Mon, 15 Sep 2014 23:49:02 (-40.6d)
3	11	nhix5haze ‡ *	2000\$	3.78682	254	Mon, 15 Sep 2014 16:50:01 (-76.3d)
4	†38	ChoKo Team 🕸		3.77526	216	Mon, 15 Sep 2014 15:21:36 (-42.1h)
5	†35	cheng chen		3.77384	21	Mon, 15 Sep 2014 23:29:29 (-0h)
6	†16	quantify		3.77086	8	Mon, 15 Sep 2014 16:12:48 (-7.3h)
7	11	Stanislav Seme	nov & Co (HSE Yandex)	3.76211	68	Mon, 15 Sep 2014 20:19:03
8	47	Luboš Motl's te	am 🕫	3.76050	589	Mon, 15 Sep 2014 08:38:49 (-1.6h)
9	†8	Roberto-UCIIIM		3.75864	292	Mon, 15 Sep 2014 23:44:42 (-44d)
10	† <b>2</b>	Davut & Josef 🗐		3.75838	161	Mon, 15 Sep 2014 23:24:32 (-4.5d)
45	15	crowwork ± ‡	HEP meets ML award Free trip to CERN	3.71885	94	Mon, 15 Sep 2014 23:45:00 (-5.1d)
782	Ļ149	Eckhard	TMVA expert, with TM improvements	VA 3.4994	5 29	Mon, 15 Sep 2014 07:26:13 (-46.1h)
991	14	Rem.		3.20423	2	Mon, 16 Jun 2014 21:53:43 (-30.4h)
.9.						e. 26 Jan 2022 15

Yann Coadou (CPPM) — Machine learning

## 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' -lTMVA
- 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");

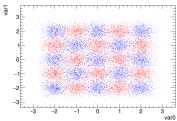


TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!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);

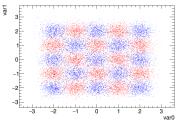




TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!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 = "";

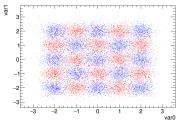


TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!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 =
                                                       0È
   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"); TTree\* sig = (TTree\*)inputFile->Get("TreeS");

```
TTree* bkg = (TTree*)inputFile->Get("TreeB");
double sigWeight = 1.0; double bkgWeight = 1.0;
```

TMVA::DataLoader \*dataloader =

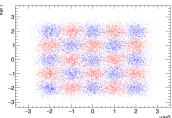
new TMVA::DataLoader("dataset"); dataloader->AddSignalTree(sig, sigWeight); dataloader->AddBackgroundTree(bkg, bkgWeight); dataloader->AddVariable("var0", 'F');

```
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



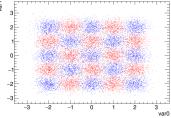


TFile\* outputFile = TFile::Open("output.root","RECREATE"); TMVA::Factory \*factory = new TMVA::Factory( "TMVAClassification", outputFile,

"!V:Color:DrawProgressBar:Transformations=I: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');



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

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");
factory->BookMethod(dataloader_TMVA::Types::kBDT___"H:IV:
```

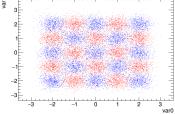
factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400: MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80"); factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher"); factory->TrainAllMethods(); // Train MVAs using training events factory->TestAllMethods(); // 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"); for the local state of the lo



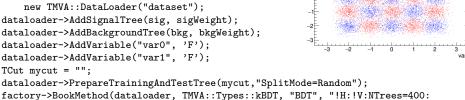
factory->BookMethod(dataloader, TMVA::Types::kBDT, "BDT", "!H:!V:NTrees=400: MinNodeSize=4%:MaxDepth=5:BoostType=AdaBoost:AdaBoostBeta=0.15:nCuts=80"); factory->BookMethod(dataloader, TMVA::Types::kFisher, "Fisher", "!H:!V:Fisher"); factory->TrainAllMethods(); // Train MVAs using training events factory->TestAllMethods(); // 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 =



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;

Yann Coadou (CPPM) — Machine learning



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<<" var1="<<var1<<" BDT="<<bdt<<" Fisher="<<fisher<<end1:</pre>
3
delete reader:
inputFile->Close();
```

More complete tutorials:

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





```
To make code compilable (and MUCH faster)
    Need ROOT and TMVA corresponding header files
    ■ e.g., for Train.C:
  #include "TFile.h"
  #include "TTree.h"
  #include "TMVA/Factory.h"
  #include "TMVA/DataLoader.h"
  #include "TMVA/TMVAGui.h"
    Need a "main" function
  int main() {
    Train():
    return 0:
  }
    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
```

