

Machine Learning and Multivariate Techniques in HEP data Analyses

Boosted Decision Trees

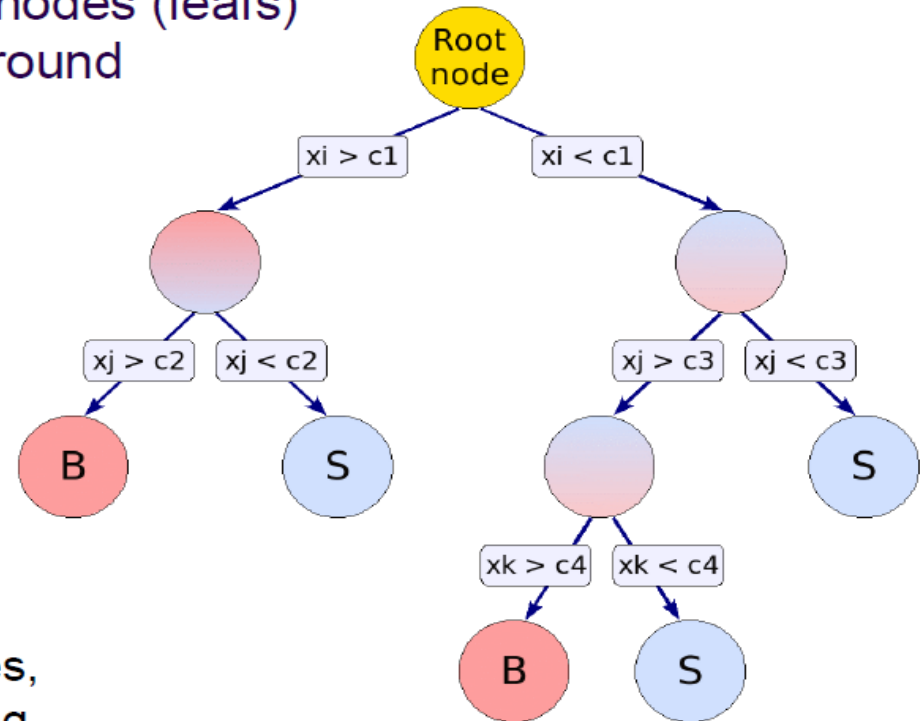
Artificial Neural Networks

Extracted from slides by:

G. Cowan's lectures at RH London Univ., H. Voss at SOS 2016, K. Reyers lectures at Heilderbeg Univ.

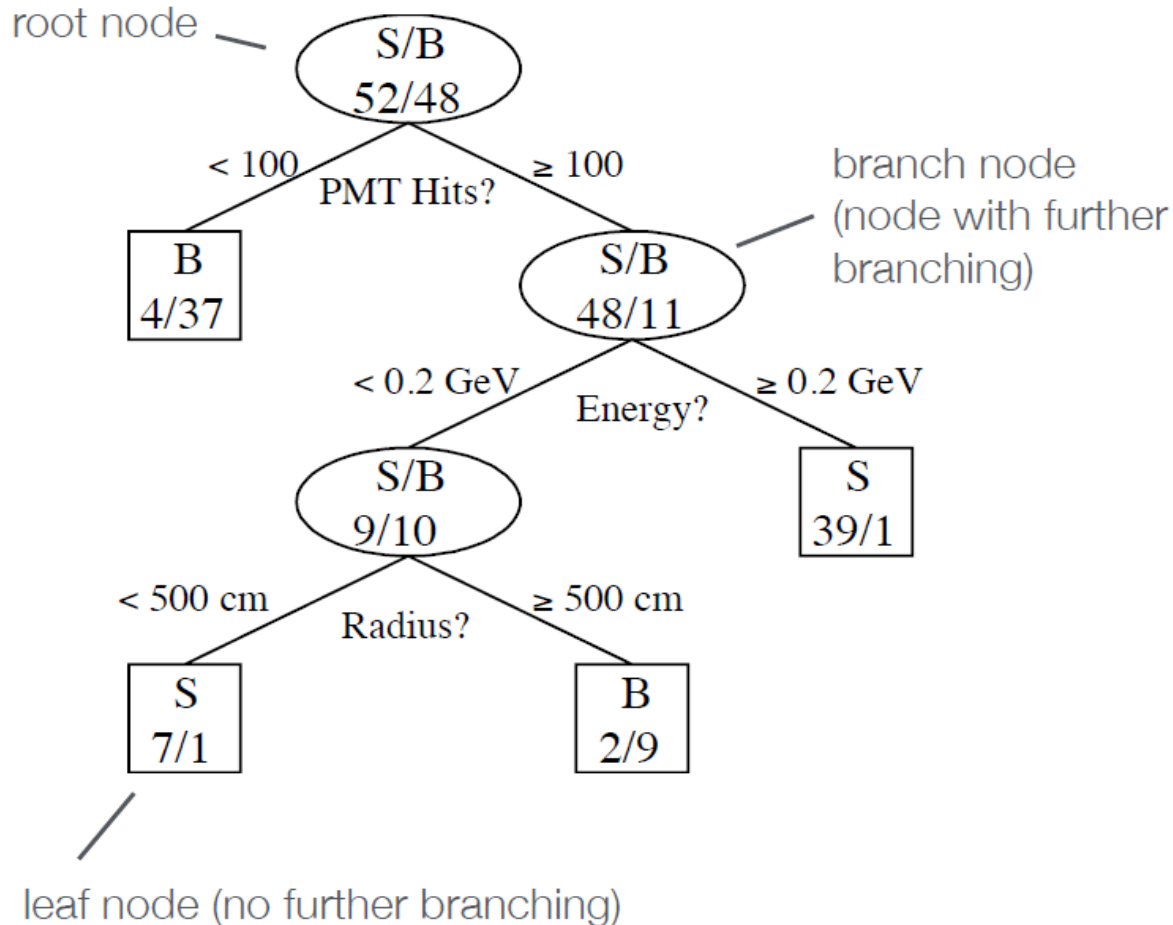
Decision Trees

- Decision Tree: Sequential application of cuts splits the data into nodes, where the final nodes (leafs) classify an event as signal or background
- Each branch \rightarrow one standard “cut” sequence
 - easy to interpret, visualised
- Disadvantage \rightarrow very sensitive to statistical fluctuations in training data
- Boosted Decision Trees (1996): combine a whole forest of Decision Trees, derived from the same sample, e.g. using different event weights.
 - overcomes the stability problem
 - increases performance

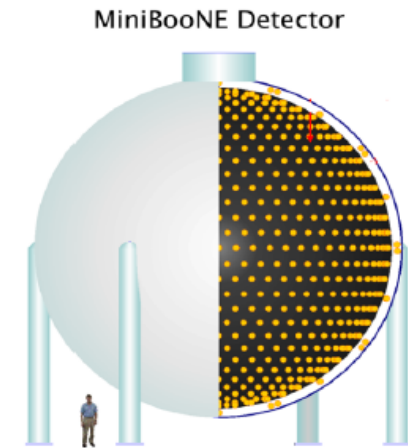


\rightarrow became popular in HEP since MiniBooNE, B.Roe et.a., NIM 543(2005)

Boosted Decision Trees



arXiv:physics/0508045v1

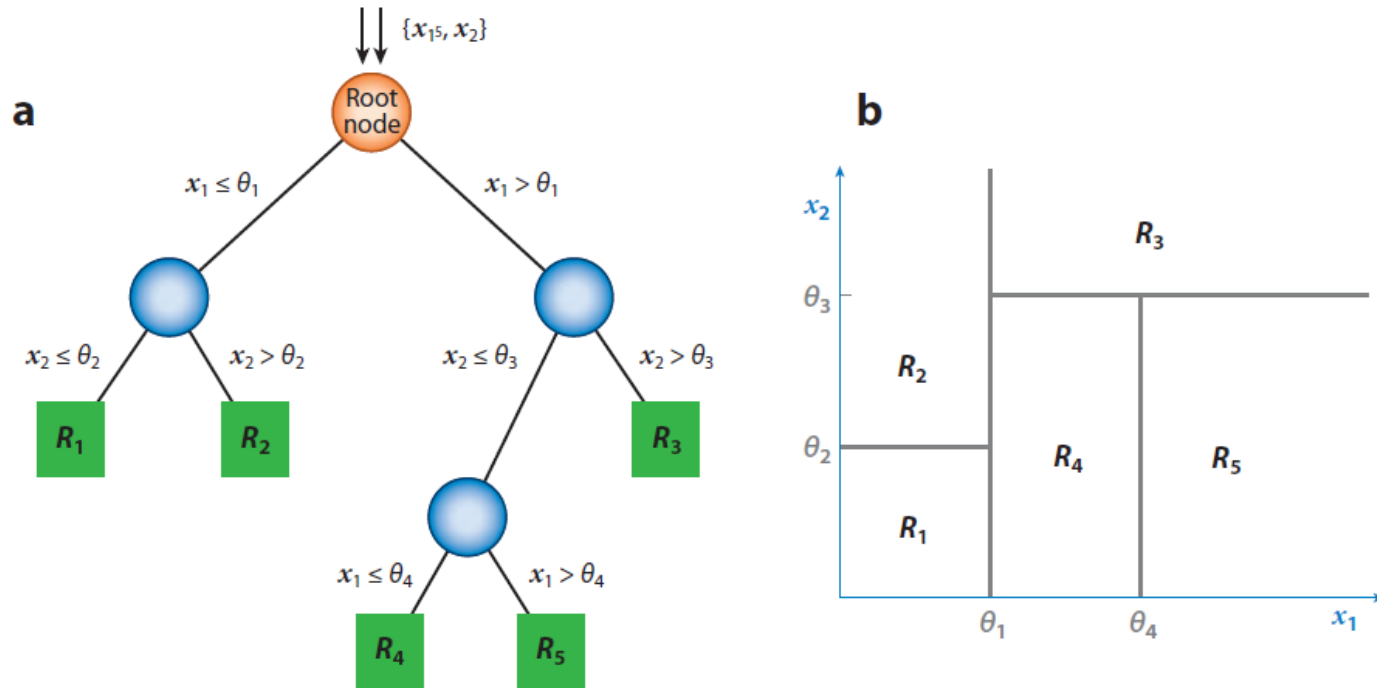


MiniBooNE: 1520 photomultiplier signals, goal: separation of ν_e from ν_μ events

Leaf nodes classify events as either signal or background

Decision Trees

Ann.Rev.Nucl.Part.Sci. 61 (2011) 281-309



Easy to interpret and visualize:

Space of feature vectors split up into rectangular volumes
(attributed to either signal or background)

How to build a decision tree in an optimal way?

Finding Optimal Cuts

Separation btw. signal and background is often measured with the *Gini index*:

$$G = p(1 - p)$$

Here p is the purity:

$$p = \frac{\sum_{\text{signal}} w_i}{\sum_{\text{signal}} w_i + \sum_{\text{background}} w_i}$$

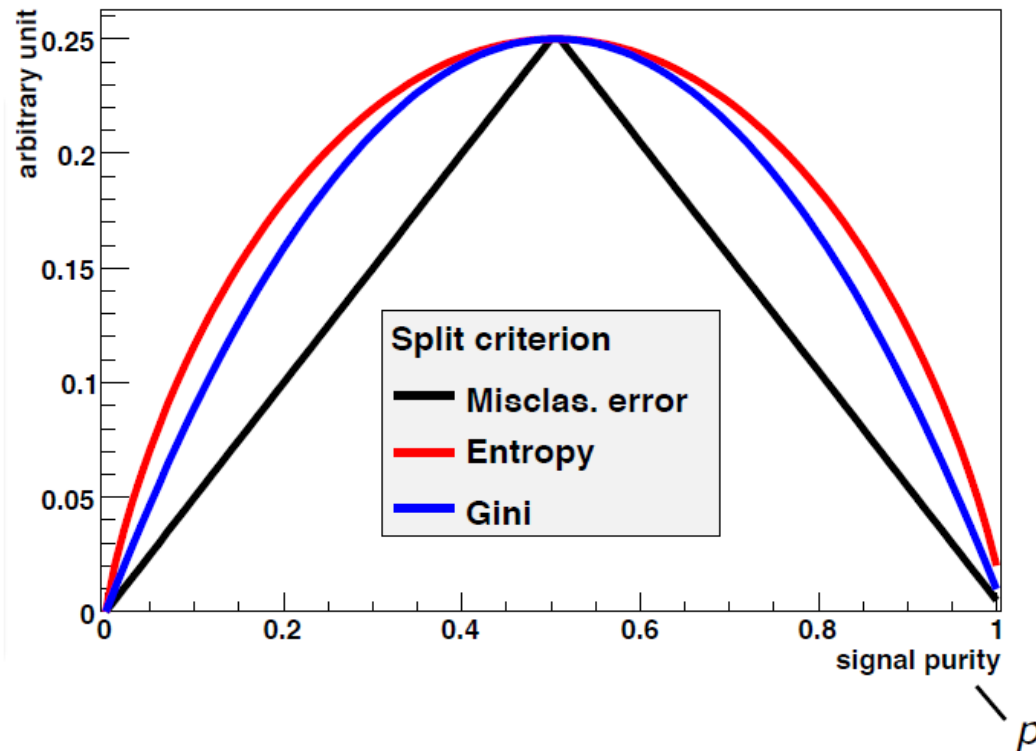
w_i = weight of event i

[usefulness of weights will become apparent soon]

Improvement in signal/background separation after splitting a set A into two sets B and C:

$$\Delta = W_A G_A - W_B G_B - W_C G_C \quad \text{where} \quad W_X = \sum_X w_i$$

Separation Measures



Cross entropy: $-p \ln p - (1 - p) \ln(1 - p)$

Gini index: $p(1 - p)$ [after Corrado Gini, used to measure income and wealth inequalities, 1912]

Misclassification rate: $1 - \max(p, 1 - p)$

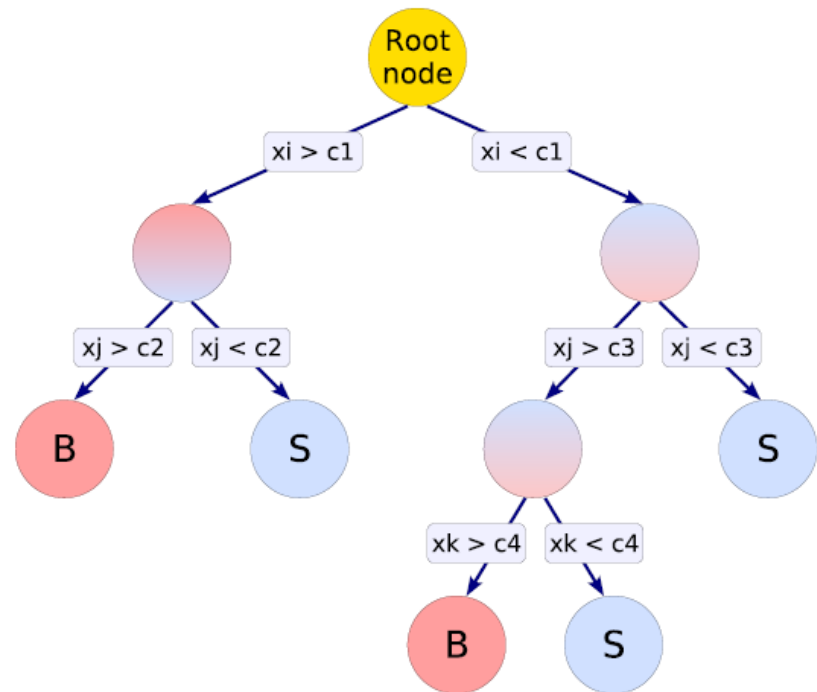
Decision Tree Pruning

When to stop growing a tree?

- ▶ When all nodes are essentially pure?
- ▶ Well, that's overfitting!

Pruning

- ▶ Cut back fully grown tree to avoid overtraining



Boosted Decision Trees: Idea

Drawback of decisions trees:

very sensitive to statistical fluctuations in training sample

Solution: boosting

- ▶ One tree → several trees ("*forrest*")
- ▶ Trees are derived from the same training ensemble by reweighting events
- ▶ Individual trees are then combined: weighted average of individual trees

Boosting is a general method of combining a set of classifiers (not necessarily decisions trees) into a new, more stable classifier with smaller error.

Popular example: AdaBoost (Freund, Schapire, 1997)

AdaBoost (short for Adaptive Boosting)

Initial training sample

$\vec{x}_1, \dots, \vec{x}_n$: multivariate event data
 y_1, \dots, y_n : true class labels, +1 or -1
 $w_1^{(1)}, \dots, w_n^{(1)}$: event weights

with equal weights normalized as

$$\sum_{i=1}^n w_i^{(1)} = 1$$

Train first classifier f_1 :

$f_1(\vec{x}_i) > 0$ classify as signal
 $f_1(\vec{x}_i) < 0$ classify as background

Assigning the Classifier Score

Assign score to each classifier according to its error rate:

$$\alpha_k = \ln \frac{1 - \varepsilon_k}{\varepsilon_k}$$

Combined classifier (weighted average):

$$f(\vec{x}) = \sum_{k=1}^K \alpha_k f_k(\vec{x})$$

It can be shown that the error rate of the combined classifier satisfies

$$\varepsilon \leq \prod_{k=1}^K 2\sqrt{\varepsilon_k(1 - \varepsilon_k)}$$

Updating Events Weights

Define training sample $k+1$ from training sample k by updating weights:

$$w_i^{(k+1)} = w_i^{(k)} \frac{e^{-\alpha_k f_k(\vec{x}_i) y_i / 2}}{Z_k}$$

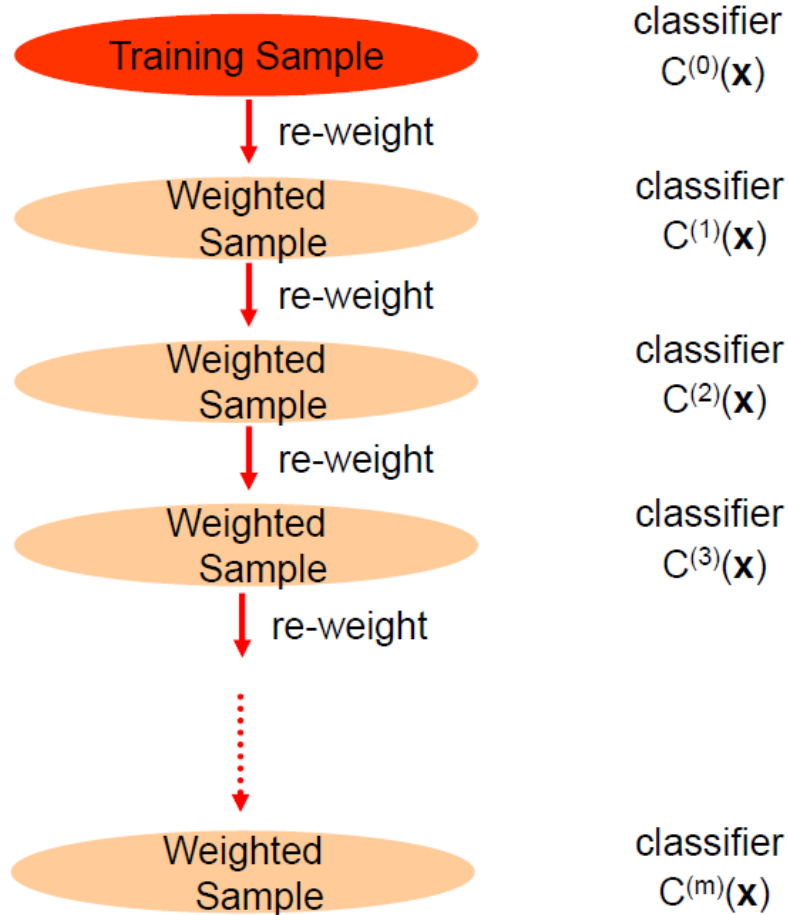
$i = \text{event index}$ normalization factor so that $\sum_{i=1}^n w_i^{(k)} = 1$

Weight is increased if event was misclassified by the previous classifier
→ "Next classifier should pay more attention to misclassified events"

At each step the classifier f_k minimizes error rate

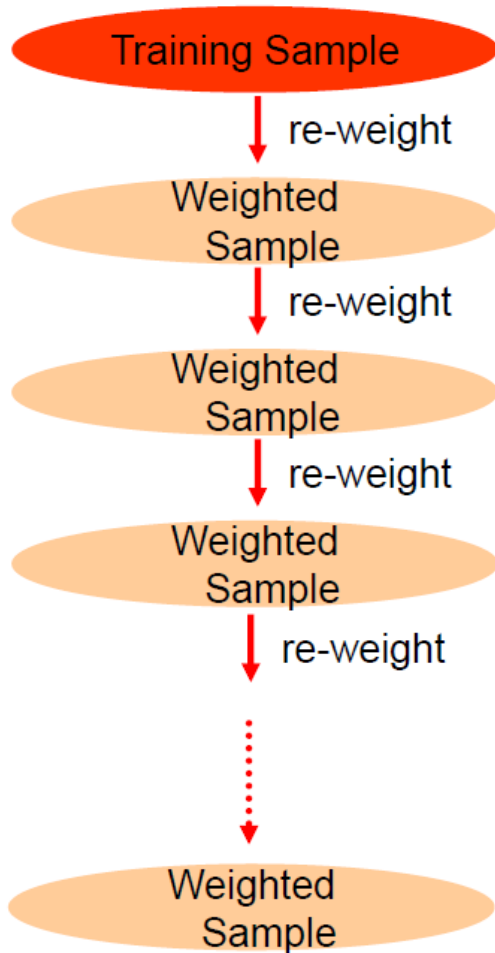
$$\varepsilon_k = \sum_{i=1}^n w_i^{(k)} I(y_i f_k(\vec{x}_i) \leq 0), \quad I(X) = 1 \text{ if } X \text{ is true, } 0 \text{ otherwise}$$

Boosting



$$y(\mathbf{x}) = \sum_i^{N_{\text{Classifier}}} w_i C^{(i)}(\mathbf{x})$$

Addaptive Boosting (AdaBoost)



classifier
 $C^{(0)}(\mathbf{x})$

classifier
 $C^{(1)}(\mathbf{x})$

classifier
 $C^{(2)}(\mathbf{x})$

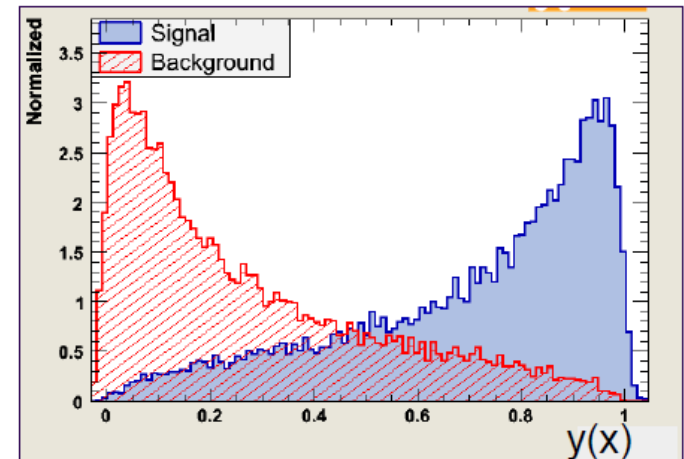
classifier
 $C^{(3)}(\mathbf{x})$

classifier
 $C^{(m)}(\mathbf{x})$

- AdaBoost re-weights events misclassified by previous classifier:

$$\frac{1 - f_{err}}{f_{err}} ; f_{err} = \frac{\text{misclassified}}{\text{all events}}$$

$$y(\mathbf{x}) = \sum_i^{N_{\text{Classifier}}} \log \left(\frac{1 - f_{err}^{(i)}}{f_{err}^{(i)}} \right) C^{(i)}(\mathbf{x})$$



Boosted Decision Trees

- Are very popular in HEP
 - Robust and easy to train,
 - get good results
- But: when we adopted BDTs,
 - In 2006 ANNs just started their big breakthrough in the ML community with remarkable advances in DEEP Learning !

General Remarks on Multi-Variate Analyses

MVA Methods

- ▶ More effective than classic cut-based analyses
- ▶ Take correlations of input variables into account

Important: find good input variables for MVA methods

- ▶ Good separation power between S and B
- ▶ Little correlations among variables
- ▶ No correlation with the parameters you try to measure in your signal sample!

Pre-processing

- ▶ Apply obvious variable transformations and let MVA method do the rest
- ▶ Make use of obvious symmetries: if e.g. a particle production process is symmetric in polar angle θ use $|\cos \theta|$ and not $\cos \theta$ as input variable
- ▶ It is generally useful to bring all input variables to a similar numerical range

Machine Learning - Basic terminology

The goal of machine learning is to predict results based on incoming data.

Features (also parameters, or variables): these are the factors for a machine to look at. E.g.: cartesian coordinates, pixel colors, a car mileage, user's gender, stock price, word frequency in the text.

- Quantitative ($x = \{1.02, 0.21, 0.12, 2\}$)
- Qualitative *discrete* ($x = \{\text{medium, small, large}\}$) or *categorical* ($x = \{\text{red, blue, green}\}$)

Algorithms (also models): Any problem can be solved in different ways. The method you choose affects the precision, performance, and size of the final model.

- If the data is insufficient/inappropriate (e.g. statistically limited or missing important features), even the best algorithm won't help. Pay attention to the accuracy of your results only when you have a good enough dataset.

CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

SUPERVISED

Predict a category

CLASSIFICATION

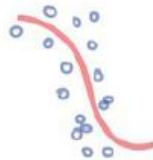
«Divide the socks by color»



Predict a number

REGRESSION

«Divide the ties by length»



OUR FOCUS

Data is not labeled in any way

UNSUPERVISED

Divide by similarity

CLUSTERING

«Split up similar clothing into stacks»



Identify sequences

Find hidden dependencies

ASSOCIATION

«Find what clothes I often wear together»



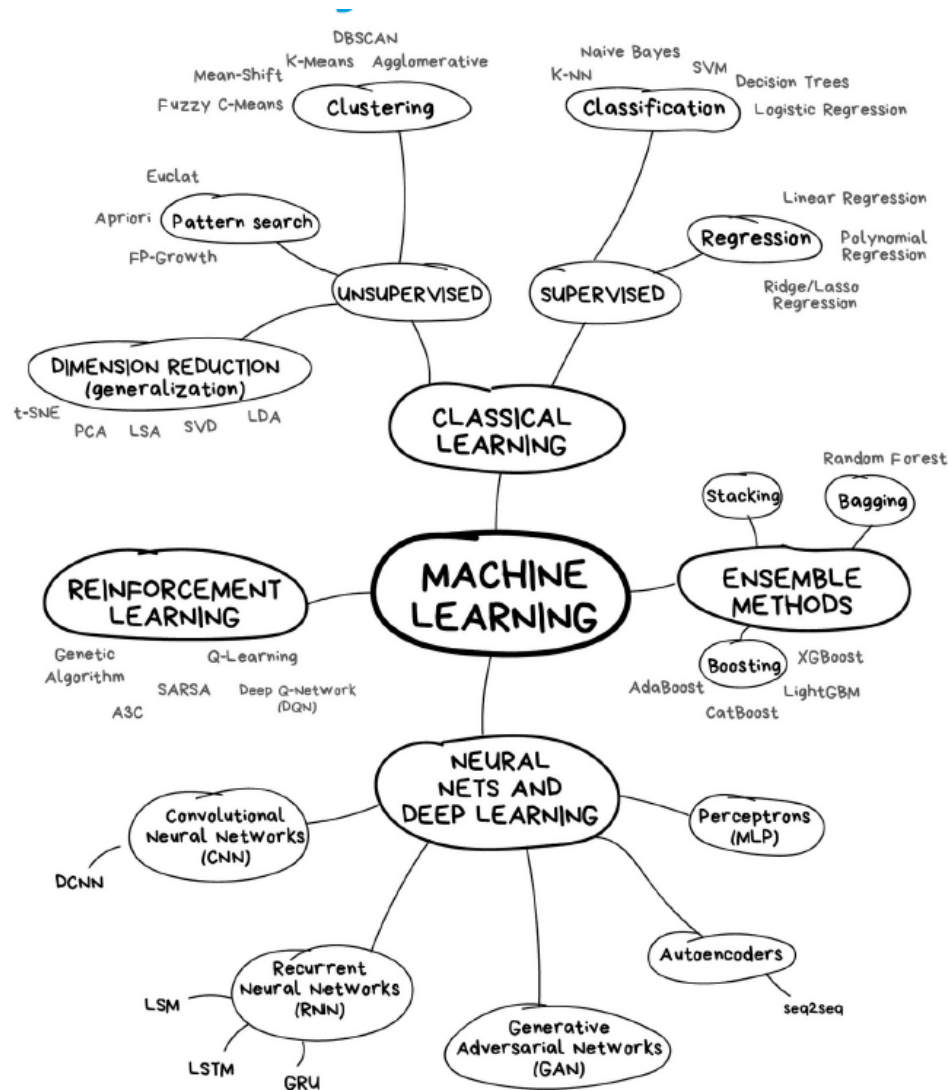
DIMENSION REDUCTION (generalization)

«Make the best outfits from the given clothes»



Image credit: https://vas3k.com/blog/machine_learning/

Where are the Neural Networks?



Neural Networks

Any neural network is a collection of **neurons** and **connections** between them.

Neuron is a function with a set of inputs and one output. Its task is to take all numbers from its input, apply a function on them and send the result to the output.

- Example: sum up all numbers from the inputs and if that sum is bigger than N give 1 as a result. Otherwise return zero.

Connections are like channels between neurons. They connect outputs of one neuron with the inputs of another so they can send digits to each other. Each connection has only one parameter the *weight*.

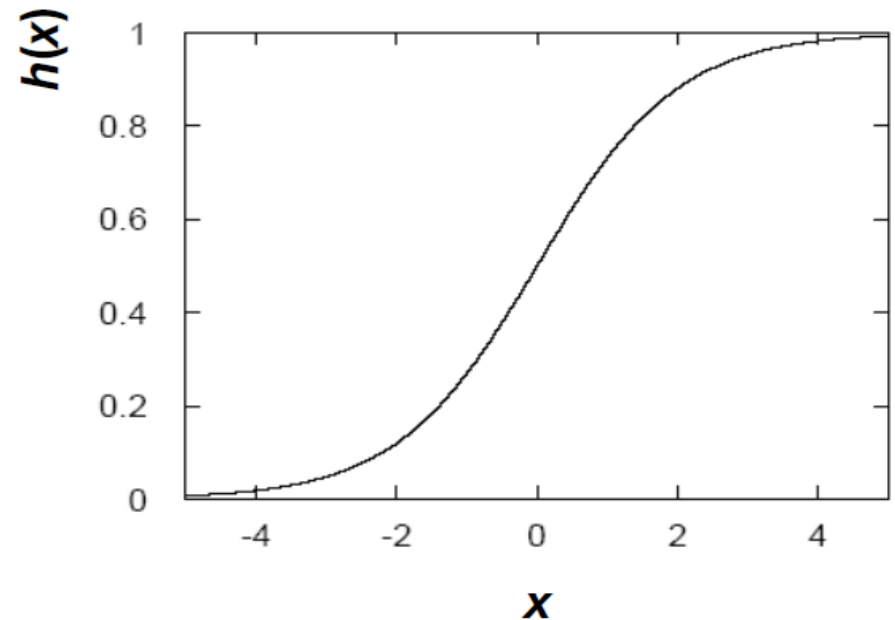
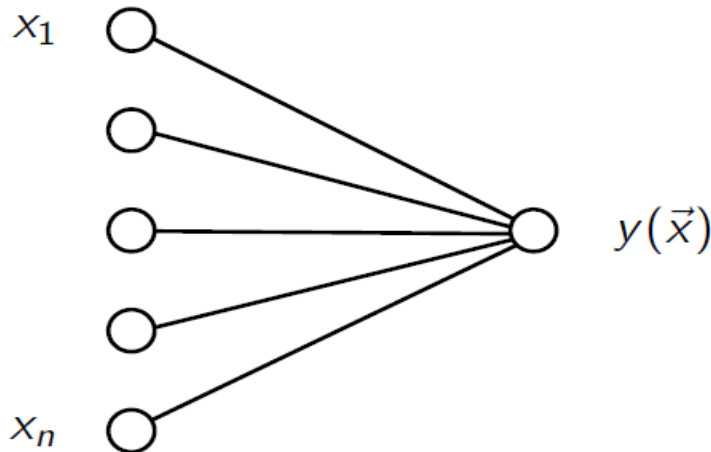
- These weights tell the neuron to respond more to one input and less to another. Weights are adjusted when training — that's how the network learns.

Perceptron

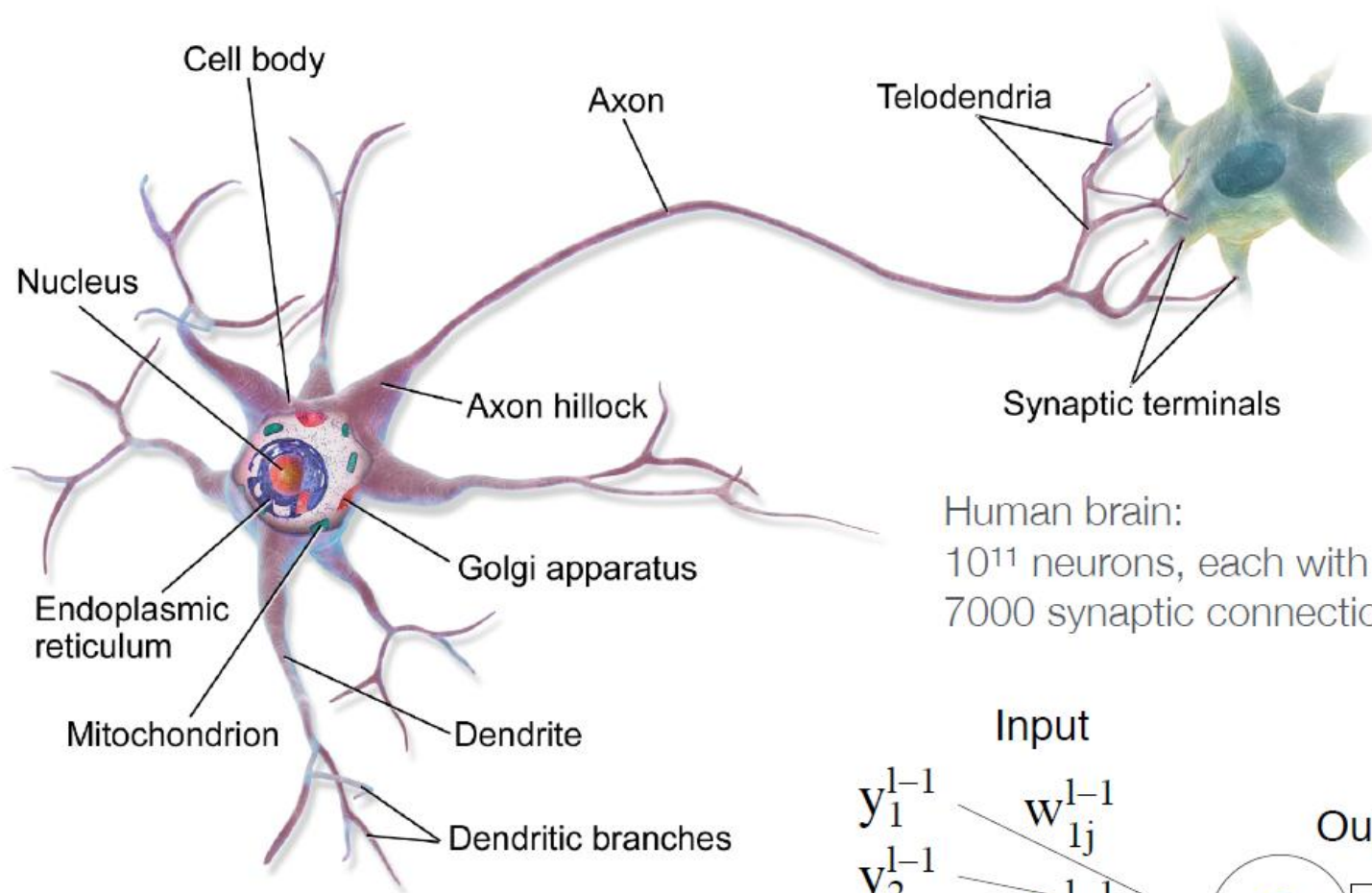
Discriminant:
$$y(\vec{x}) = h \left(w_0 + \sum_{i=1}^n w_i x_i \right)$$

The nonlinear, monotonic function h is called *activation function*.

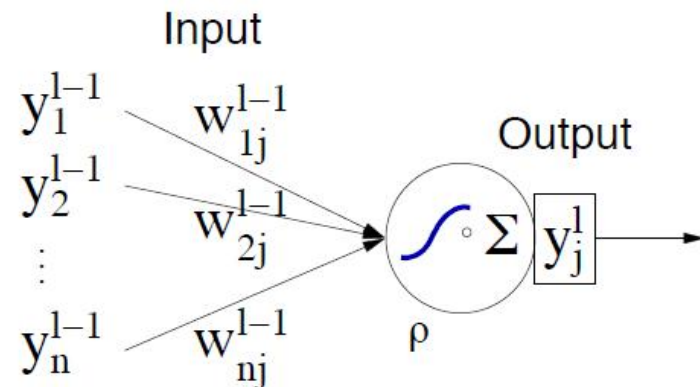
Typical choices for h : $\frac{1}{1 + e^{-x}}$ ("sigmoid"), $\tanh x$



The Biological Inspiration: the Neuron

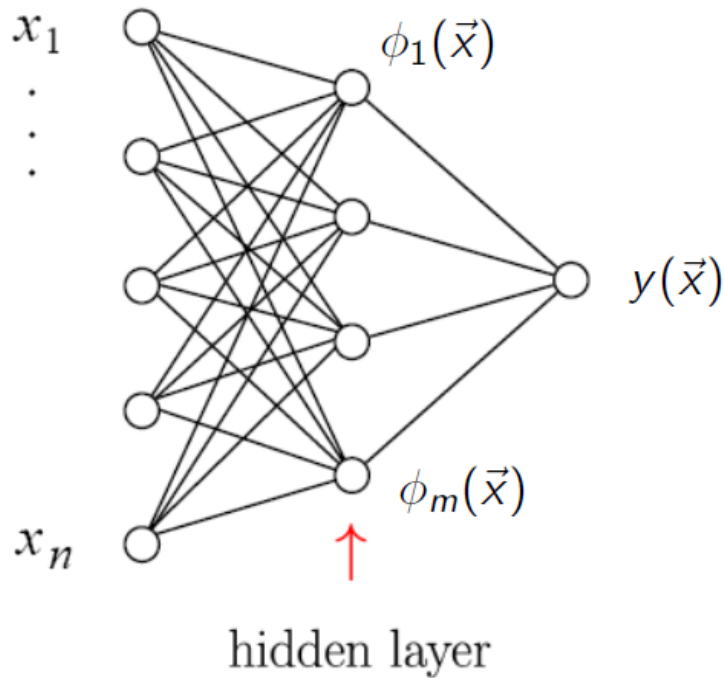


Human brain:
10¹¹ neurons, each with on average
7000 synaptic connections



<https://en.wikipedia.org/wiki/Neuron>

Feedforward Neural Network with One Hidden Layer



superscripts indicates layer number

$$\phi_i(\vec{x}) = h \left(w_{i0}^{(1)} + \sum_{j=1}^n w_{ij}^{(1)} x_j \right)$$

$$y(\vec{x}) = h \left(w_{10}^{(2)} + \sum_{j=1}^m w_{1j}^{(2)} \phi_j(\vec{x}) \right)$$

Straightforward to generalize to multiple hidden layers

Network Training

\vec{x}_a : training event, $a = 1, \dots, N$

t_a : correct label for training event a

 e.g., $t_a = 1, 0$ for signal and background, respectively

\vec{w} : vector containing all weights

Error function:

$$E(\vec{w}) = \frac{1}{2} \sum_{a=1}^N (y(\vec{x}_a, \vec{w}) - t_a)^2 = \sum_{a=1}^N E_a(\vec{w})$$

Weights are determined by minimizing the error function.

Backpropagation

Start with an initial guess $\vec{w}^{(0)}$ for the weights and then update weights after each training event:

$$\vec{w}^{(\tau+1)} = \vec{w}^{(\tau)} - \eta \nabla E_a(\vec{w}^{(\tau)})$$

└── learning rate

Let's write network output as follows:

$$y(\vec{x}) = h(u(\vec{x})) \quad \text{with} \quad u(\vec{x}) = \sum_{j=0}^m w_{1j}^{(2)} \phi_j(\vec{x}), \quad \phi_j(\vec{x}) = h\left(\sum_{k=0}^n w_{jk}^{(1)} x_k\right) \equiv h(v_j(\vec{x}))$$

Here we defined $\phi_0 = x_0 = 1$ and the sums start from 0 to include the offsets.

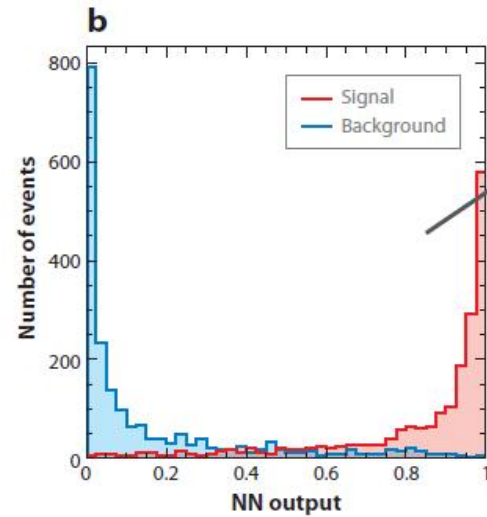
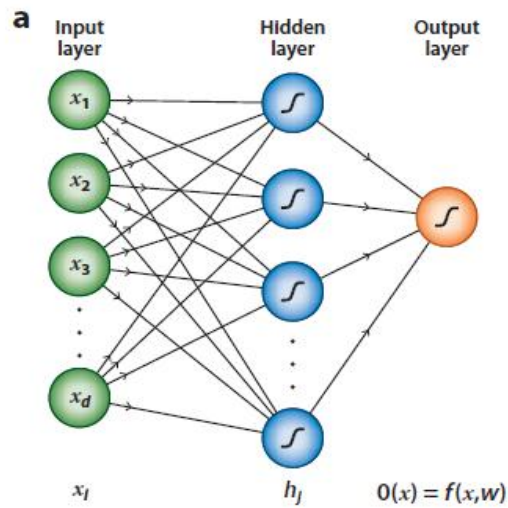
Weights from hidden layer to output:

$$E_a = \frac{1}{2}(y_a - t_a)^2 \rightarrow \frac{\partial E_a}{\partial w_{1j}^{(2)}} = (y_a - t_a) h'(u(\vec{x}_a)) \frac{\partial u}{\partial w_{1j}^{(2)}} = (y_a - t_a) h'(u(\vec{x}_a)) \phi_j(\vec{x}_a)$$

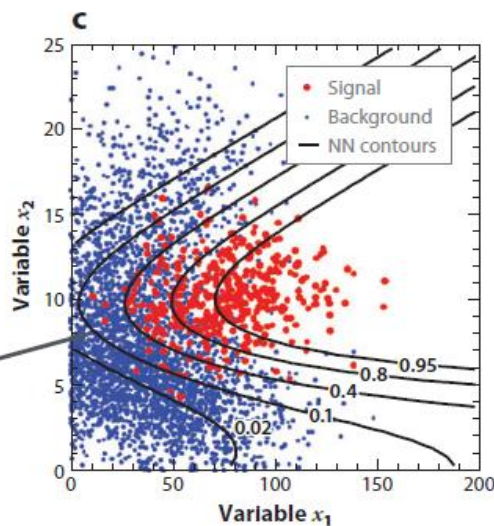
Weights from input layer to hidden layer (\rightarrow further application of chain rule):

$$\frac{\partial E_a}{\partial w_{jk}^{(1)}} = (y_a - t_a) h'(u(\vec{x}_a)) w_{1j}^{(2)} h'(v_j(\vec{x}_a)) x_{a,k} \quad \vec{x}_a \equiv (x_{a,1}, \dots, x_{a,n})$$

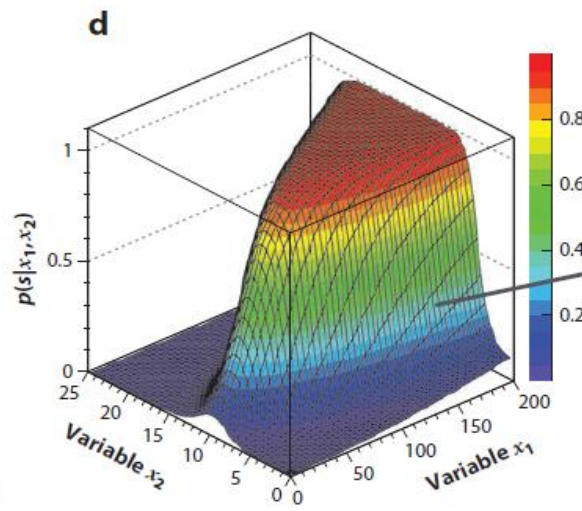
Neural Network Output and Decision Boundaries



output of neural network



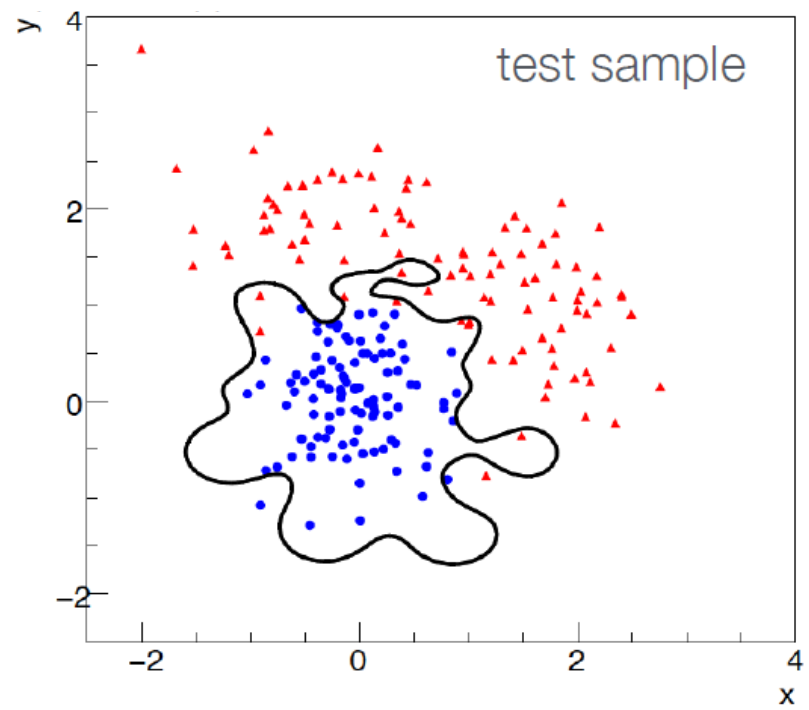
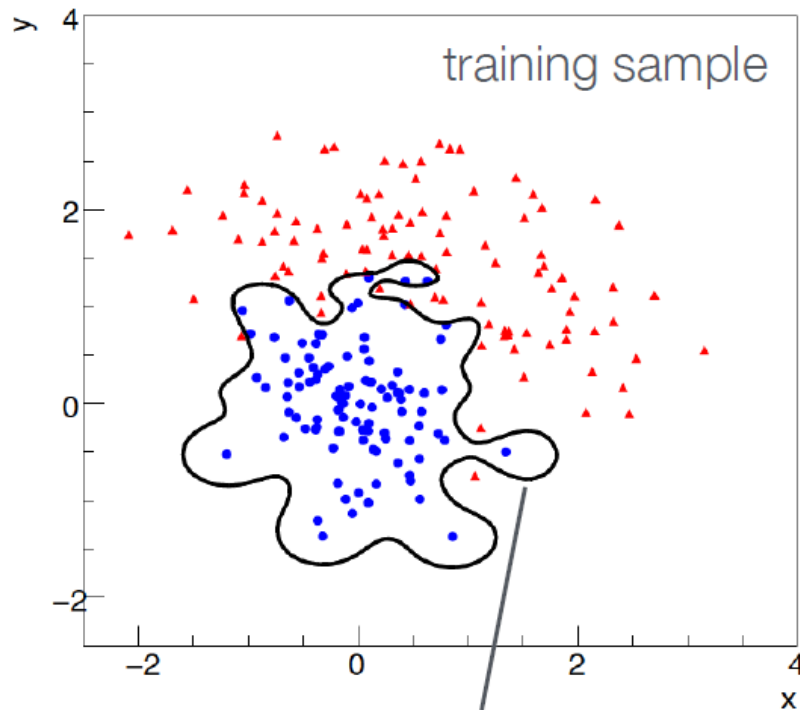
decision boundaries for different cuts on NN output



signal probability $p(s | x_1, x_2)$

Example of Overtraining

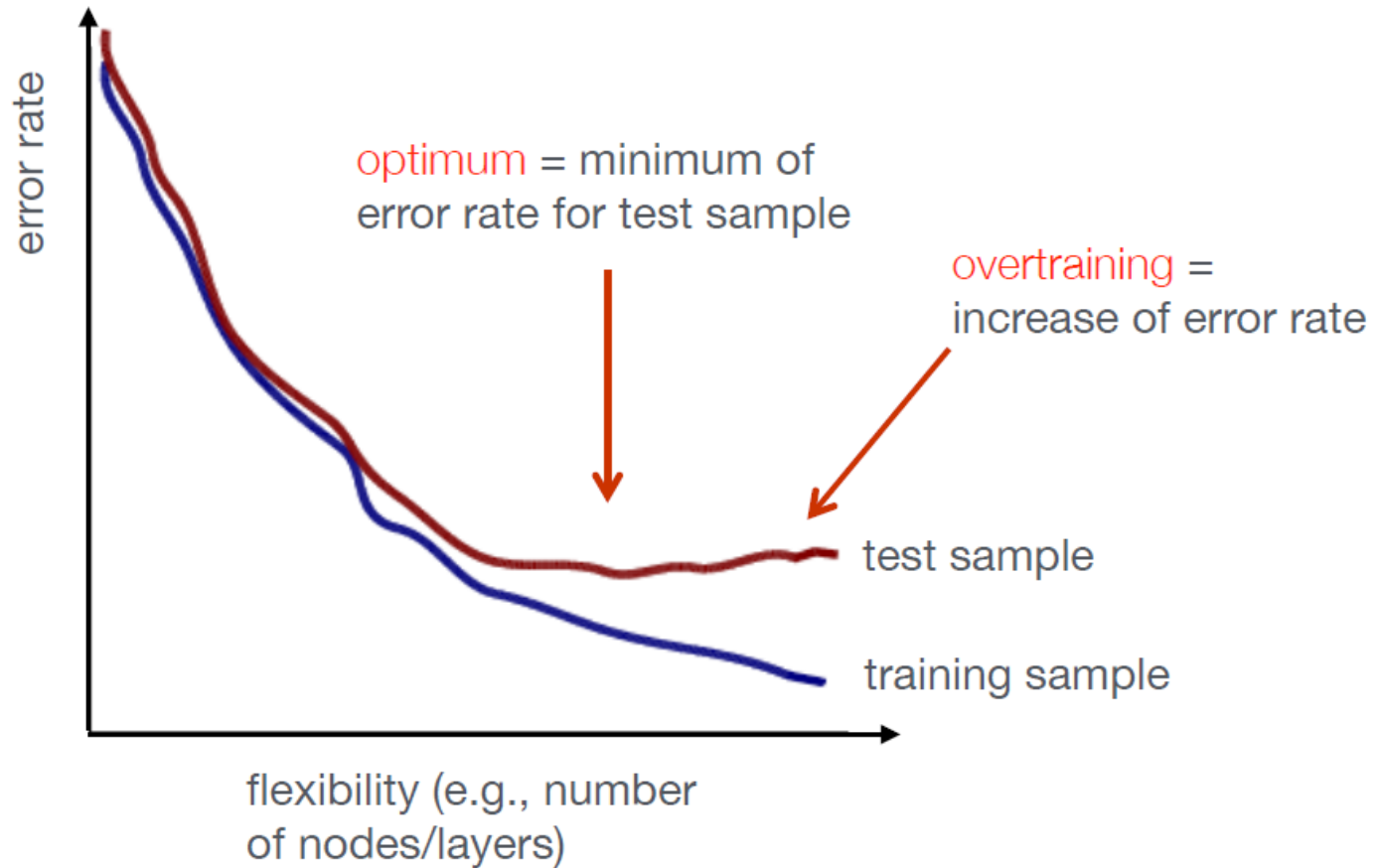
Too many neurons/layers make a neural network too flexible
→ overtraining



Network "learns" features that are merely statistical fluctuations in the training sample

Monitoring Overtraining

Monitor fraction of misclassified events (or error function:)



Deep Neural Networks

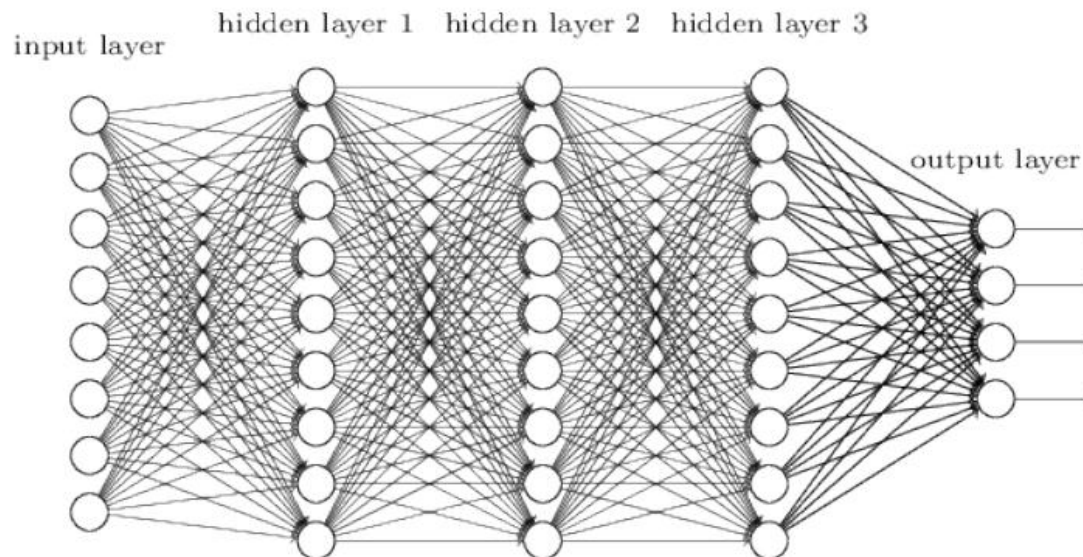
Deep networks: many hidden layers with large number of neurons

Challenges

- ▶ Hard to train ("vanishing gradient problem")
- ▶ Training slow
- ▶ Risk of overtraining

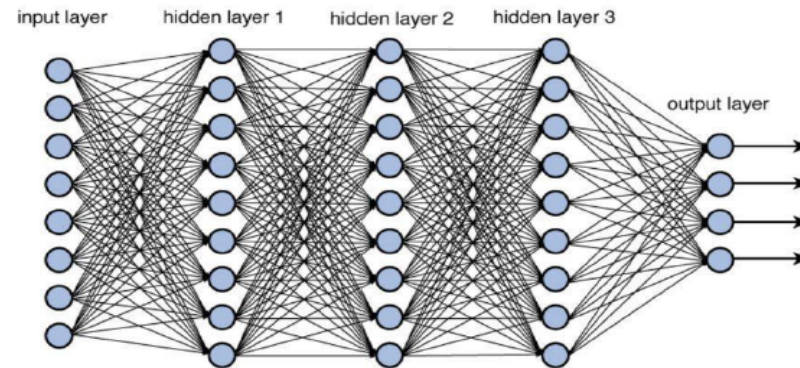
Big progress in recent years

- ▶ Interest in NN waned before ca. 2006
- ▶ Milestone: paper by G. Hinton (2006): "learning for deep belief nets"
- ▶ Image recognition, AlphaGo, ...
- ▶ Soon: self-driving cars, ...



How do NNs work?

How do NNs work?



layer

$$a_0^{(1)} = f(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \dots + w_{0,n} a_n^{(0)} + b_0)$$

activation function

weights

bias

$$\begin{bmatrix} a_0^{(1)} \\ \dots \\ a_n^{(1)} \end{bmatrix} = f \left(\begin{bmatrix} w_{0,0} & \dots & w_{0,n} \\ \vdots & \ddots & \vdots \\ w_{k,0} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ \dots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ \dots \\ b_n \end{bmatrix} \right)$$

How do NNs learn?

After we constructed a network, our task is to **assign proper weights** so neurons will react correctly to incoming signals.

- define a loss function to measure how far the response is from the truth

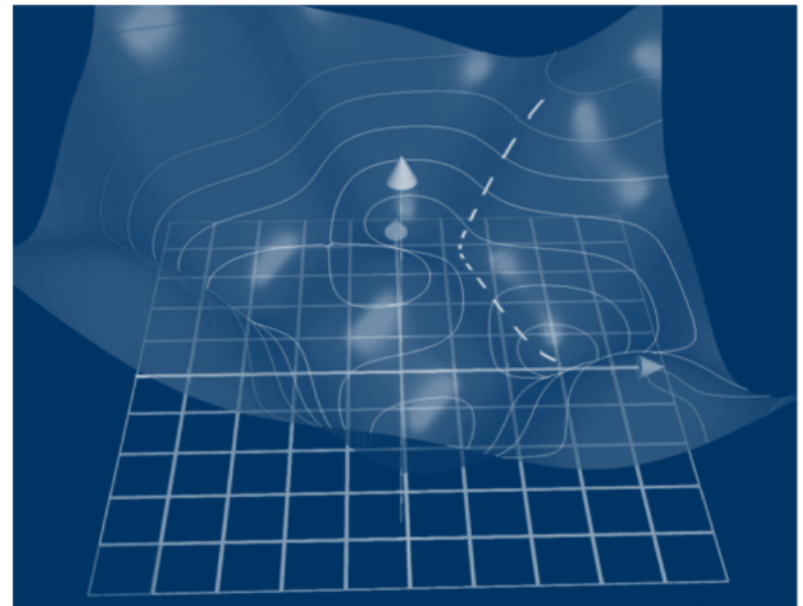
This function is a function of all the weights and biases in the NN (a priori a very large number), and the goal of training is to find its minimum.

- To start with, all weights are assigned randomly.

- After evaluating the NN on the training dataset, we can compute all the per-neuron differences with respect to the correct result.

- Computing the gradient of the loss, gives us a direction in which to tune the weights towards a local minimum

The process of correcting the weights is called backpropagation an error.



How do NNs learn?

A mostly complete chart of Neural Networks

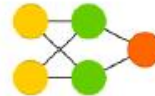
©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

- Input Cell
- Backfed Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Gated Memory Cell
- Kernel
- Convolution or Pool

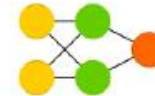
Perceptron (P)



Feed Forward (FF)



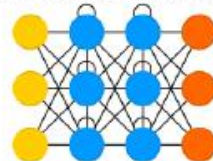
Radial Basis Network (RBF)



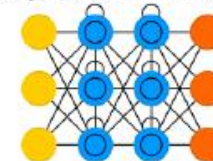
Deep Feed Forward (DFF)



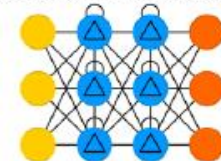
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



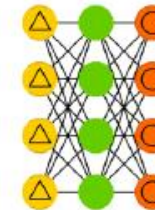
Auto Encoder (AE)



Variational AE (VAE)



Denosing AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)



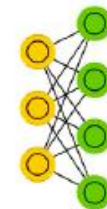
Hopfield Network (HN)



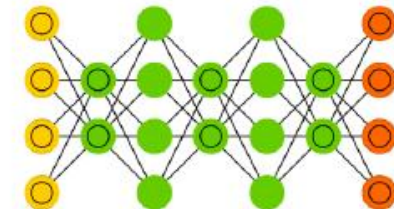
Boltzmann Machine (BM)



Restricted BM (RBM)



Deep Belief Network (DBN)



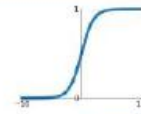
There are many more...

How do NNs learn?

- ◆ Each input is multiplied by a weight.
- ◆ The weighted values are summed and a bias is added
- ◆ **The result is passed to an activation function (non linearity).**

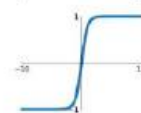
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

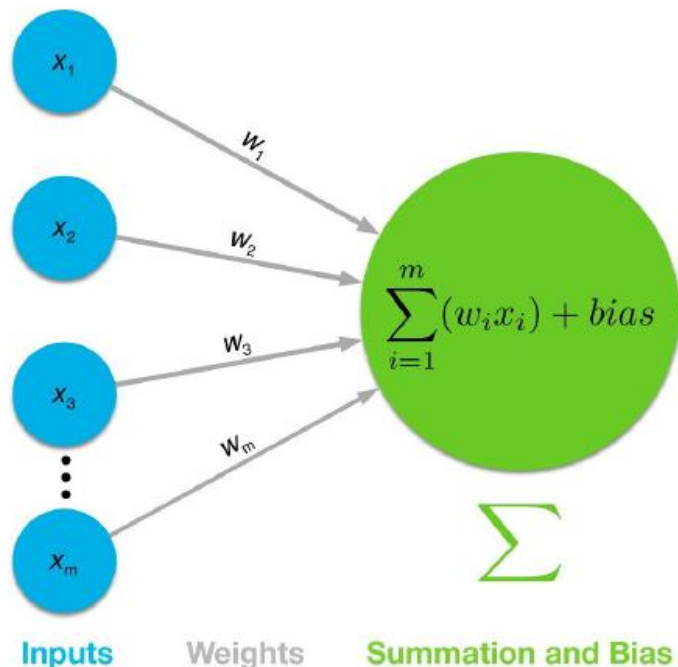
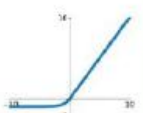


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



$$\hat{y} = f \left(\sum_{i=1}^m (w_i x_i) + bias \right)$$

$$f(x) = \begin{cases} 1 & \text{if } \sum wx + b \geq 0 \\ 0 & \text{if } \sum wx + b < 0 \end{cases}$$

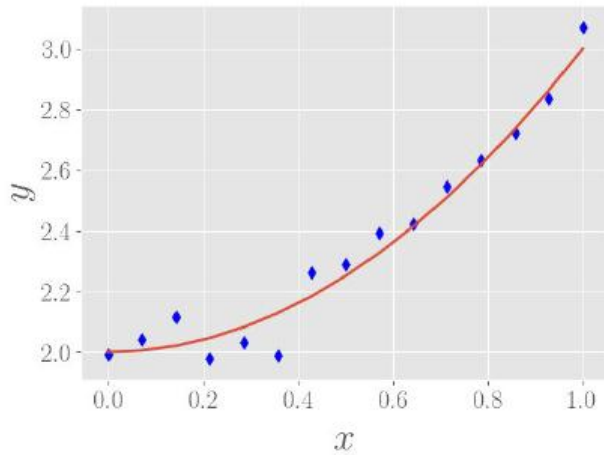
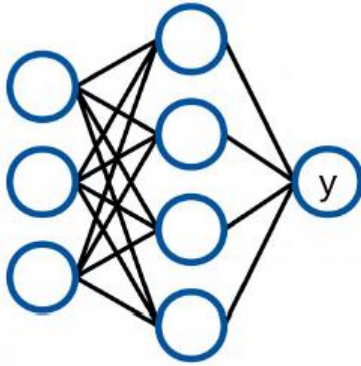


There's a colorful world inside black boxes!

Typical Applications

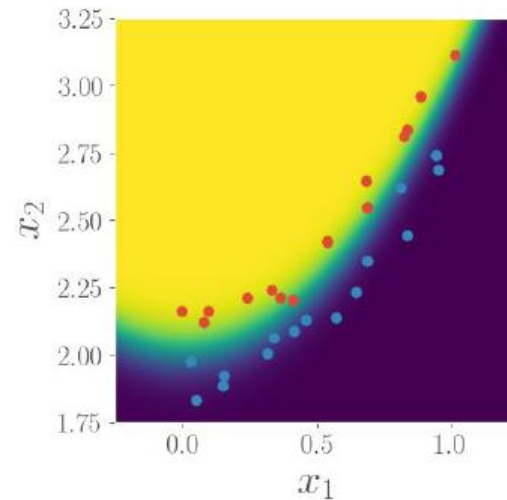
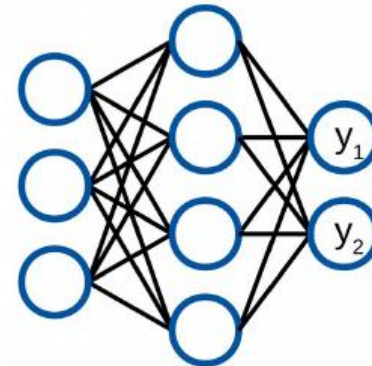
Regression:

Predict a continuous label.



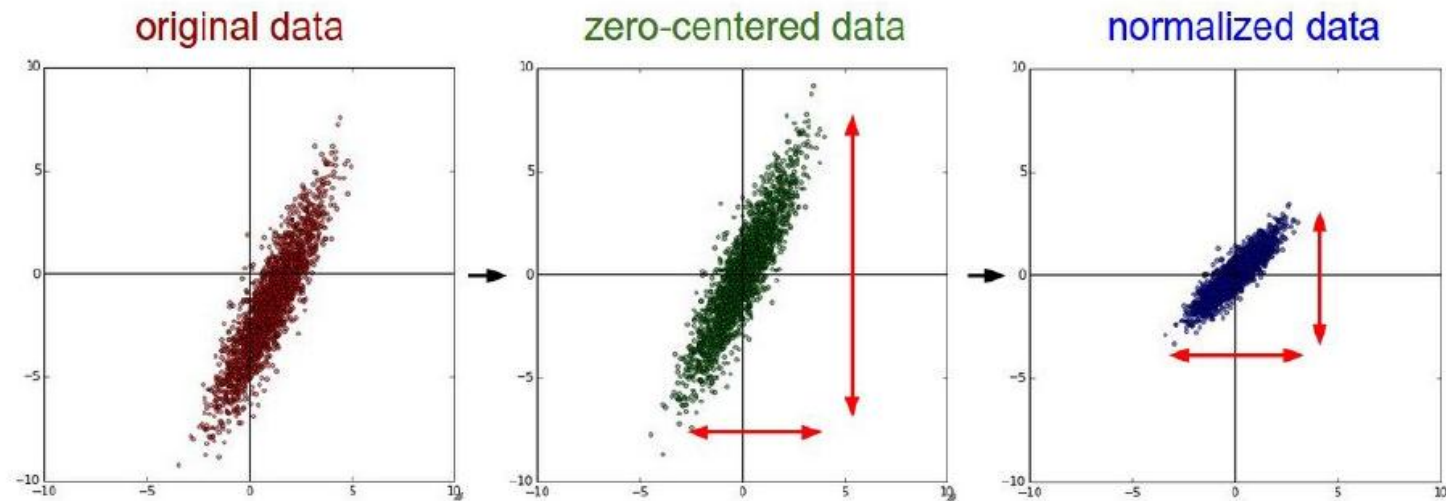
Classification:

Separate events into multiple categories.



Input Preprocessing

- ◆ Input features could have vastly different scales (e.g. p_T vs η).
 - ◆ Difficult to find optimal values.
- ◆ Basic strategy: Normalize to mean=0 and variance=1.

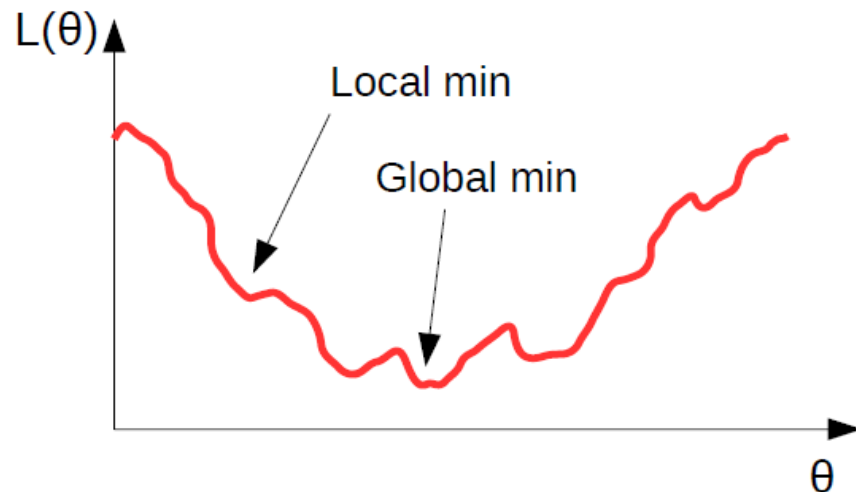


- ◆ Other options possible: decorrelation, non-linear transformation, ...

Training

- ◆ Training starts specifying an input and a target dataset.
 - ◆ For each input set, the target is what the network should learn for that input.
- ◆ A loss function is required $L(\theta)$:
 - ◆ The loss function quantifies the mistakes the NN makes. E.g. *mean squared error* can be used for regression.

Training is the minimization of the loss function w.r.t. the NN parameters.



Training: (Stochastic) Gradient Descent

- Given the increasing size of datasets and parameters, it is no more possible to directly minimize the loss function.

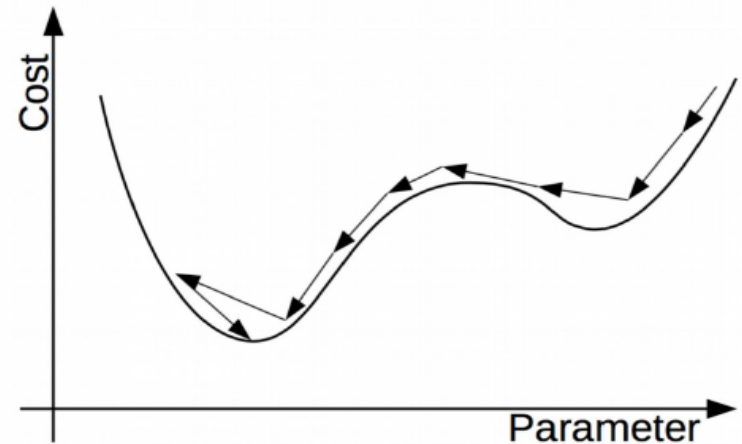
Iterative minimization by updating θ in opposite direction of gradient.

$$\theta_i = \theta_{i-1} - \alpha \frac{\partial L}{\partial \theta}, \text{ where } \alpha \text{ is the so-called « learning rate ».}$$

- Evaluation and derivation of the loss function for the full dataset is costly:

Stochastic gradient descent:

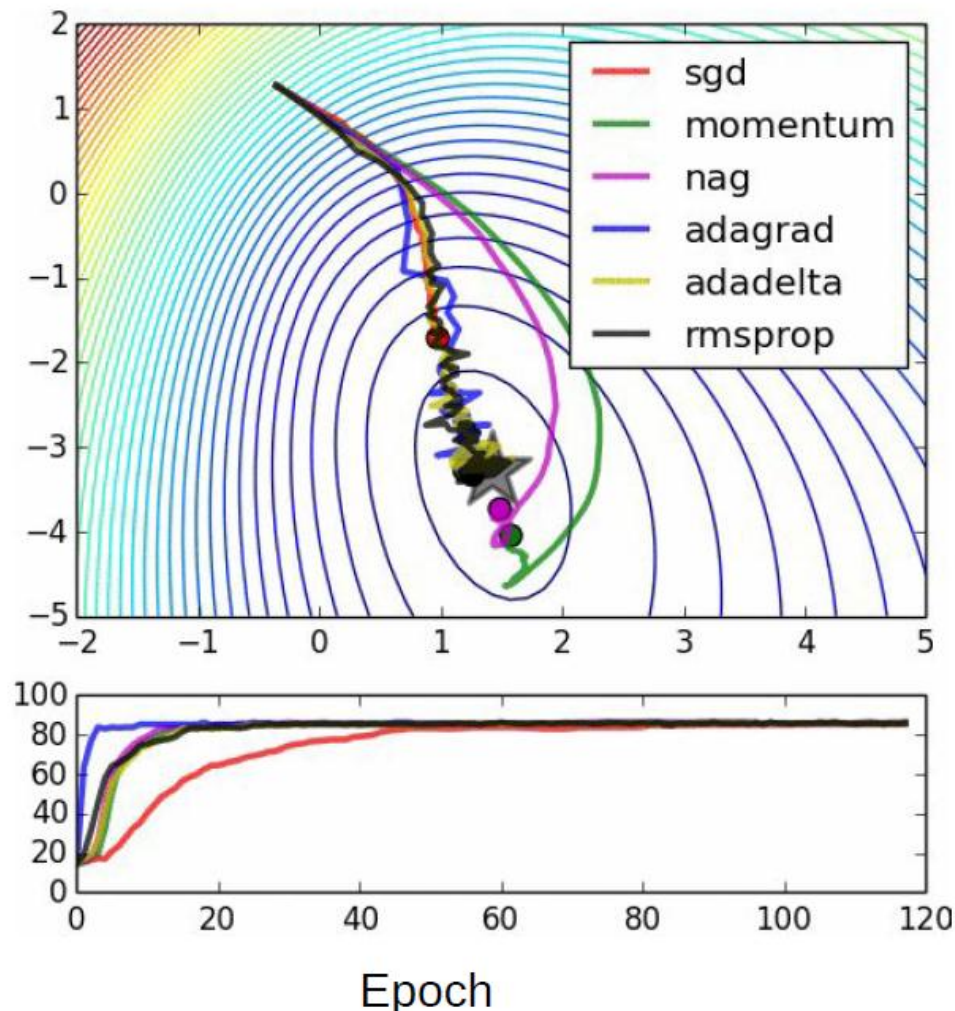
Calculate gradient for a small stochastic subset of the training sample (batch). → This also helps to avoid local minima!



One iteration over the full training dataset is called *epoch*.

Training: more optimisers

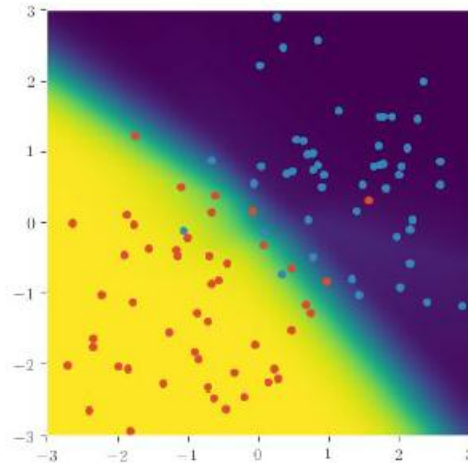
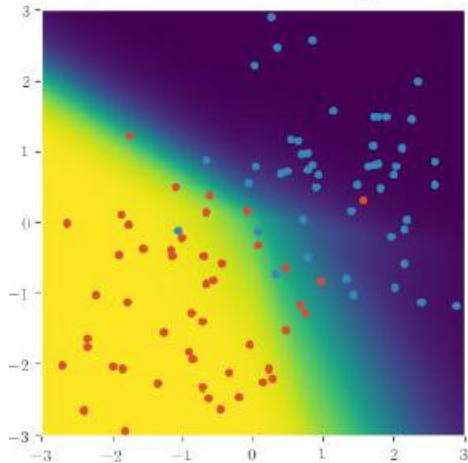
- ▶ More advanced options than fixed learning rate.
 - ◆ Momentum: past gradients used as « velocity »
 - ◆ Adaptive methods: different learning rates for each parameter and as a function of past gradients.



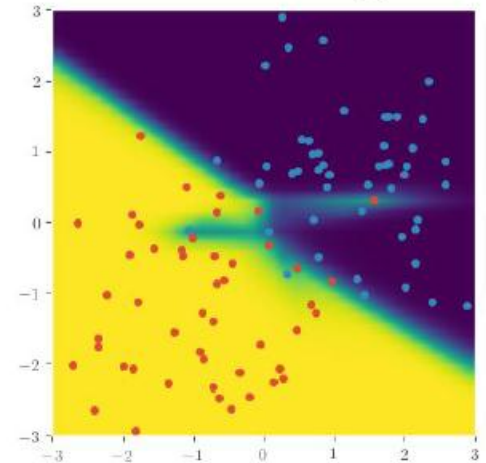
Underfitting and overtraining

Underfitting: If model capacity is too low or if training is not enough
→ bad performance.

Underfitting



Overfitting



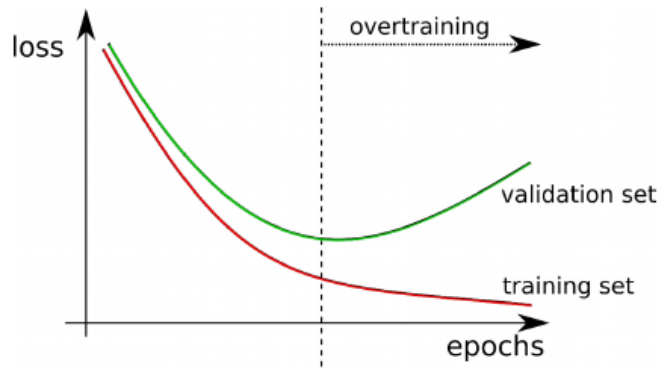
Overfitting: If model capacity is too high, network can « memorize » training samples
→ bad generalization.

Overtraining solutions

Early stop:

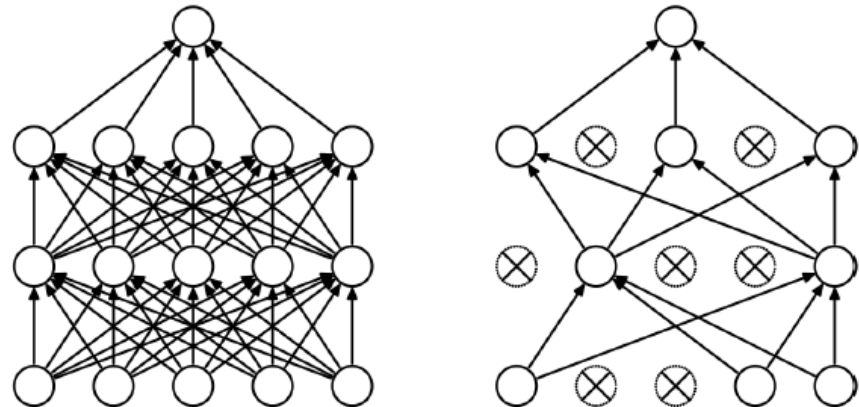
Evaluate the performance of the network on a validation dataset.

Stop when performance on validation set decreases



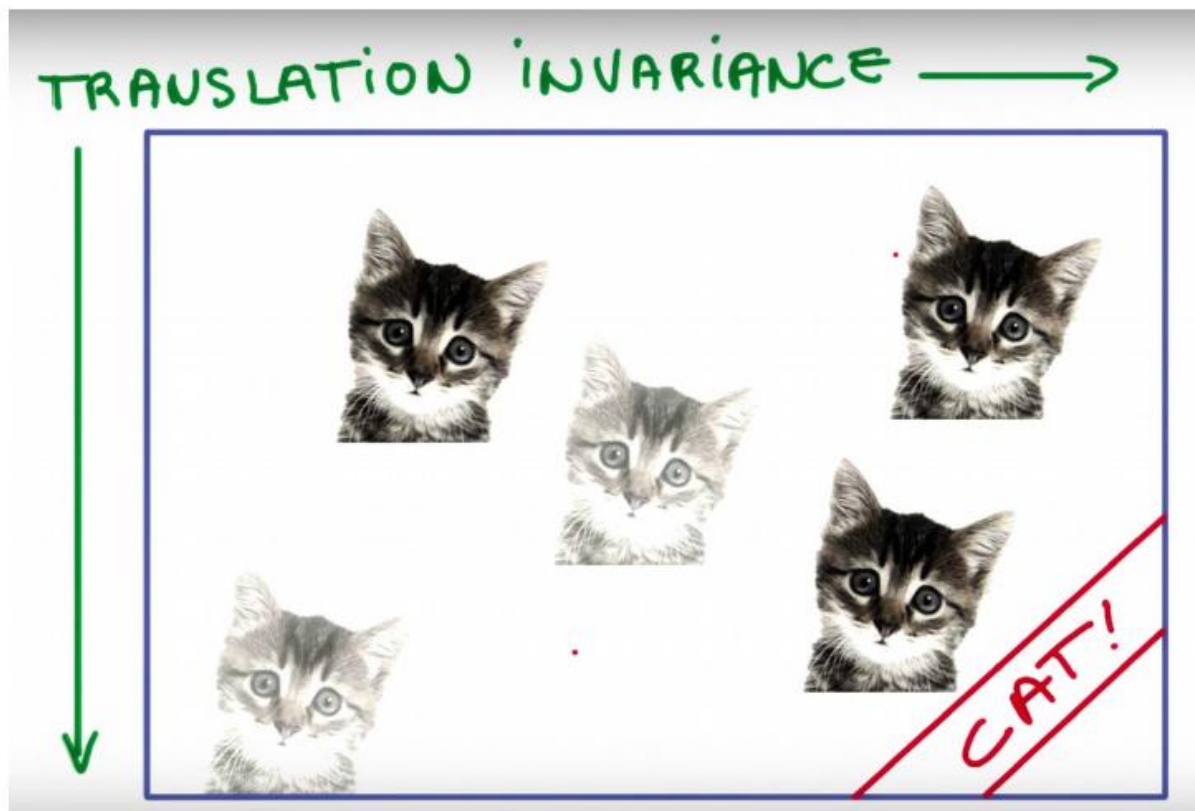
Dropout:

Randomly drop a percentage of nodes at each training step. Learn redundant representations, hence giving a more robust model.



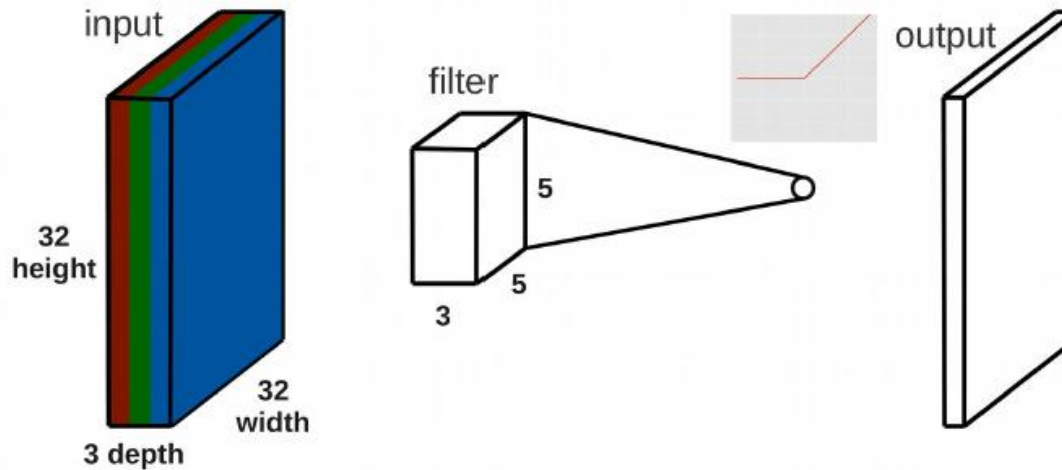
Convolutional NN

- ◆ Convolutional NNs are made to exploit local correlation and translation invariance.
- ◆ Typical applications are image processing and computer vision.



Convolution layer

- ◆ A small filter (weight tensor) slides over the image to create a feature map.



- ◆ Several filters could be stacked depth-wise.
- ◆ Several convolution can be applied one after the other to extract higher level features.

Average and Max pooling layers

- ◆ When output size reduction is required:
 - ◆ Max pooling: takes the maximum of each patch.
 - ◆ Average pooling: takes the average of each patch.
- ◆ Eg: 2x2 filter with stride=2

3	2	1	0
0	5	3	0
9	4	3	1
2	1	3	1

max pooling



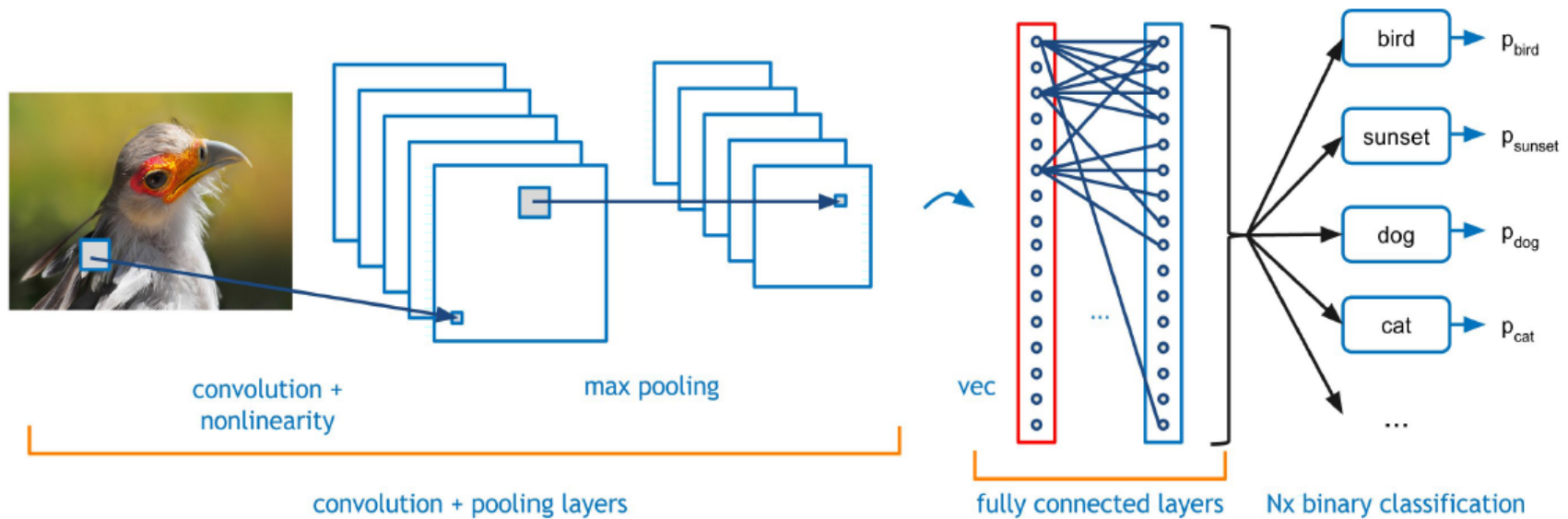
5	3
9	3

average pooling



2.5	1
4	2

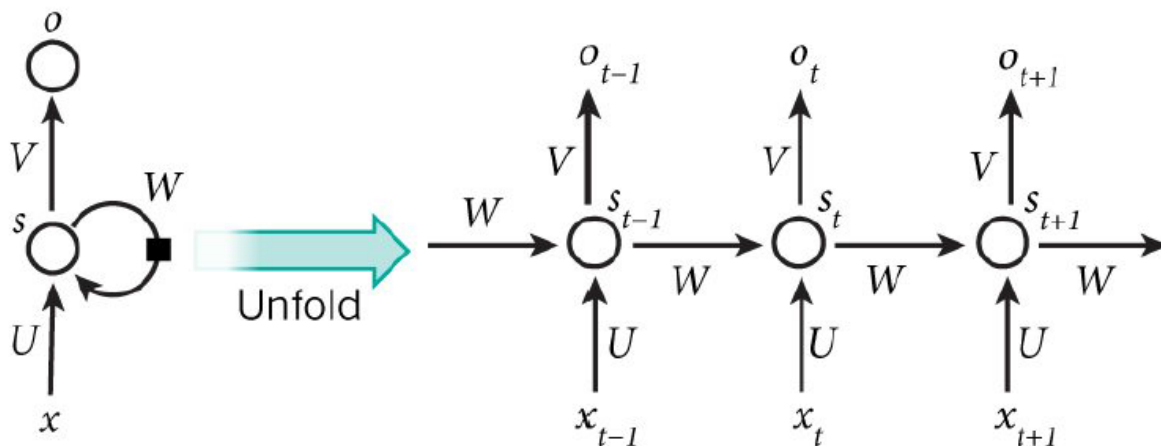
Convolutional NN architecture



- ◆ Convolution and pooling layers to extract features.
- ◆ Fully connected layers used at the end to combine features.
- ◆ Applications in HEP: PID for neutrino experiments, jet tagging, reduction of seeds for tracking, etc..

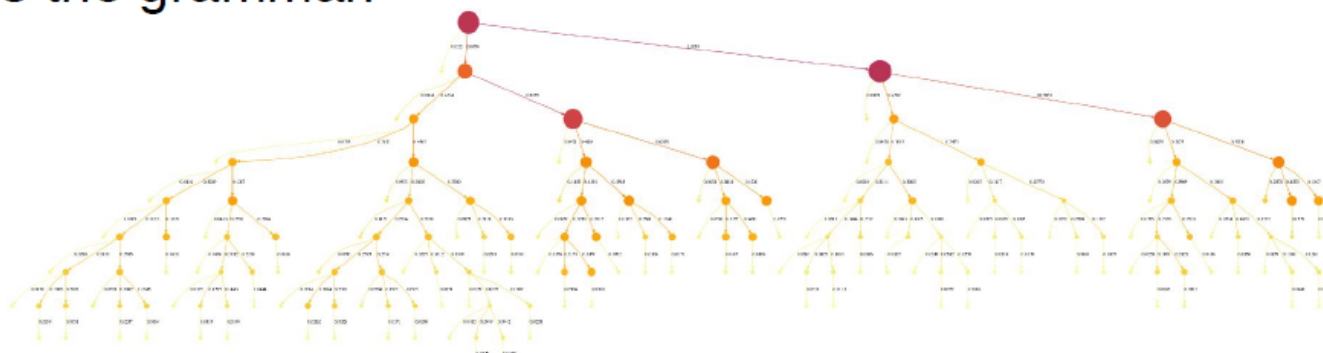
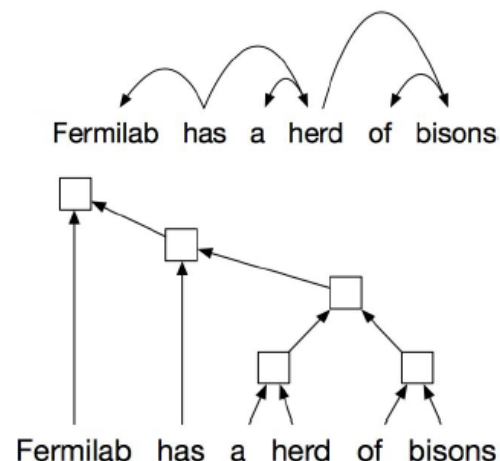
Recursive NN

- ◆ Recursive NNs are deep NN created by applying the same set of weights recursively over a structured input of variable size.
- ◆ They are called recurrent because they perform the same task for every element of a sequence, with the output being depended on the previous computations.
- ◆ Typical applications in natural language processing: apply the recursive NN to each word in a sentence for text generation (predict the next word in the sequence), translation, etc..



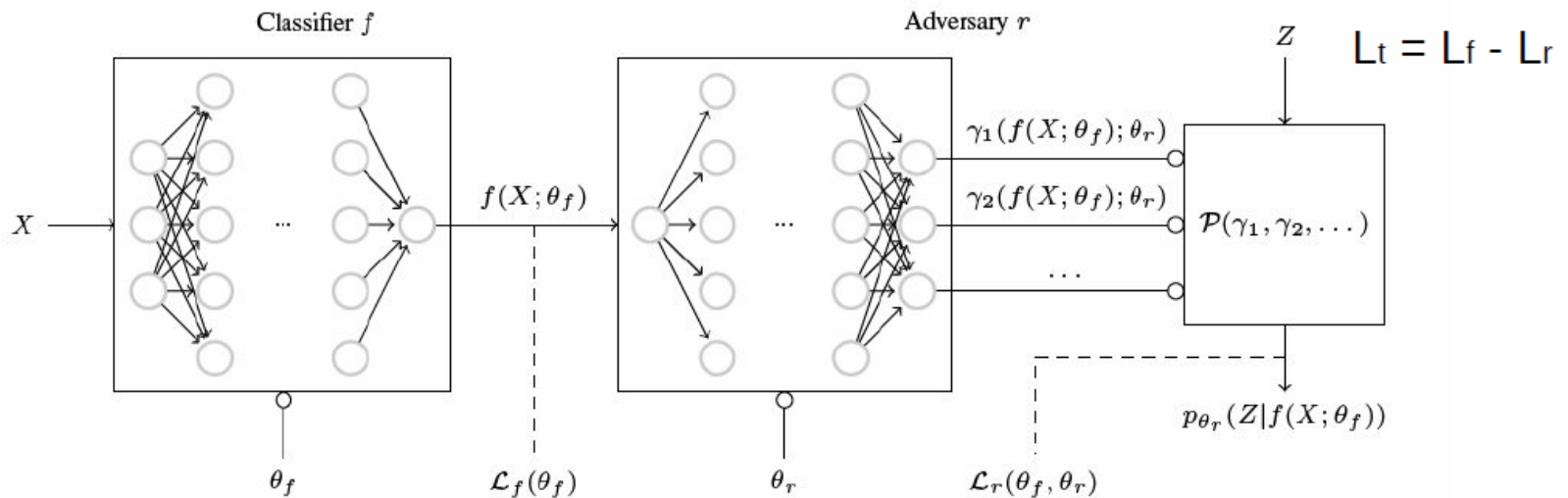
Recursive NN: possible HEP applications

- ◆ With particle-flow, collision raw data is converted in a list of particles.
- ◆ Complex objects (e.g. jets) are reconstructed by combining particles from this list.
- ◆ Image-processing approaches might not be the best in this case.
- ◆ Recursive NNs can be better suited.
 - ◆ Particles are like words in a sentence.
 - ◆ QCD is the grammar.



Adversarial NN

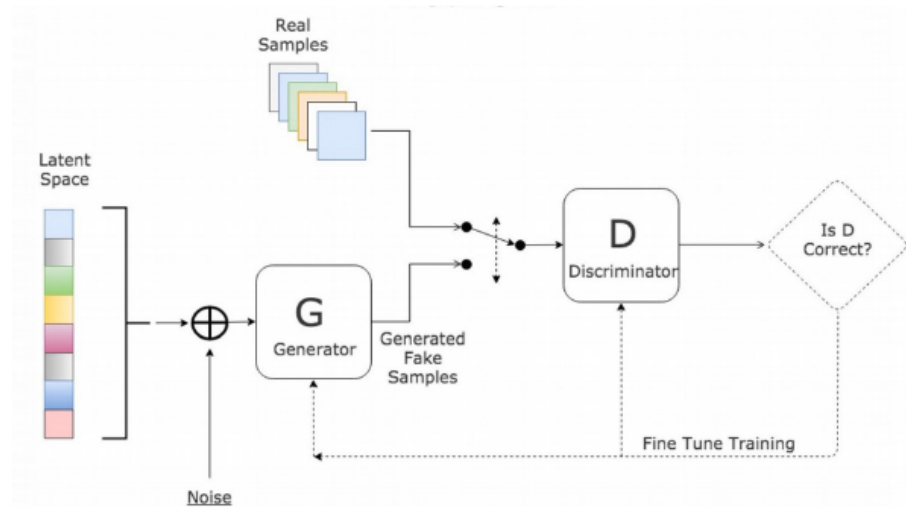
- ◆ Two deep NNs in competition with each other.
- ◆ The first NN can be used to maximize the classification performance of signal against background events.
- ◆ The second NN can be trained to identify dependency on systematic uncertainty of the output of the first NN.
- ◆ The minimization of the global loss function guarantees optimal classification performance with reduced systematic dependence.



Generative Adversarial NN

- ◆ Adversarial NNs can also be used for image generation.
- ◆ A generator and a discriminator are trained in competition.

- ◆ **Generator:**
creates images starting from noise.
- ◆ **Discriminator:**
tries to distinguish between true and generated images.



- ◆ The global loss function is given by $\text{Loss}(\text{gen}) - \text{Loss}(\text{discr})$.
- ◆ **The generator learns to make images that it has never seen simply by fooling the discriminator!**
- ◆ Application: gen of calo images, jets, and even high-level variables!



Lorentz boost network: motivation

- ◆ Deep learning methods using high-level and low-level variables are outperforming shallow learning methods using high-level variables.
 - ◆ There is some information in low-level variables that high-level variables is not using.
- ◆ Deep learning methods using low-level variables only are not able to reproduce the same results as deep learning methods using also high-level variables.
 - ◆ Need of a new NN architecture to fully exploit low-level information and automatize the design of high-level variables.

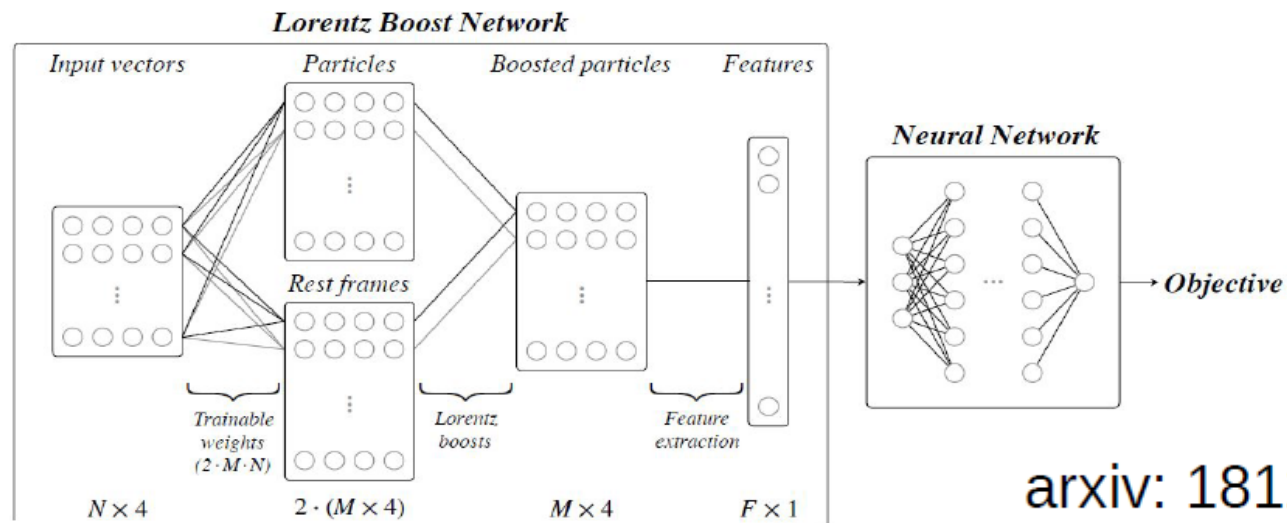
arxiv: 1812.09722

Lorentz boost network: network architecture

- ◆ Two stages approach:

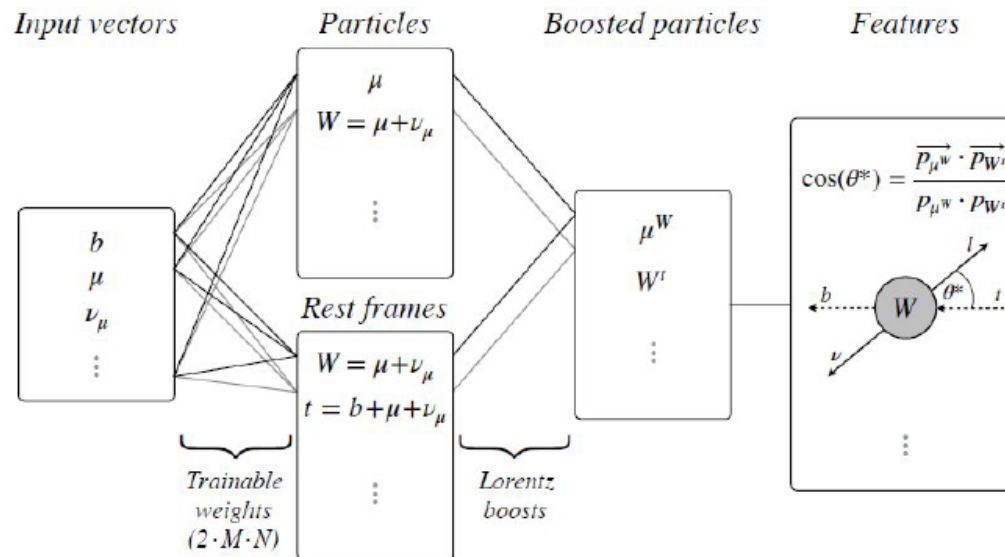
combines them to form composite particles and rest frames. Composite particles are boosted in the rest frames where features are extracted.

- ◆ An application specific NN uses LBN features as input.



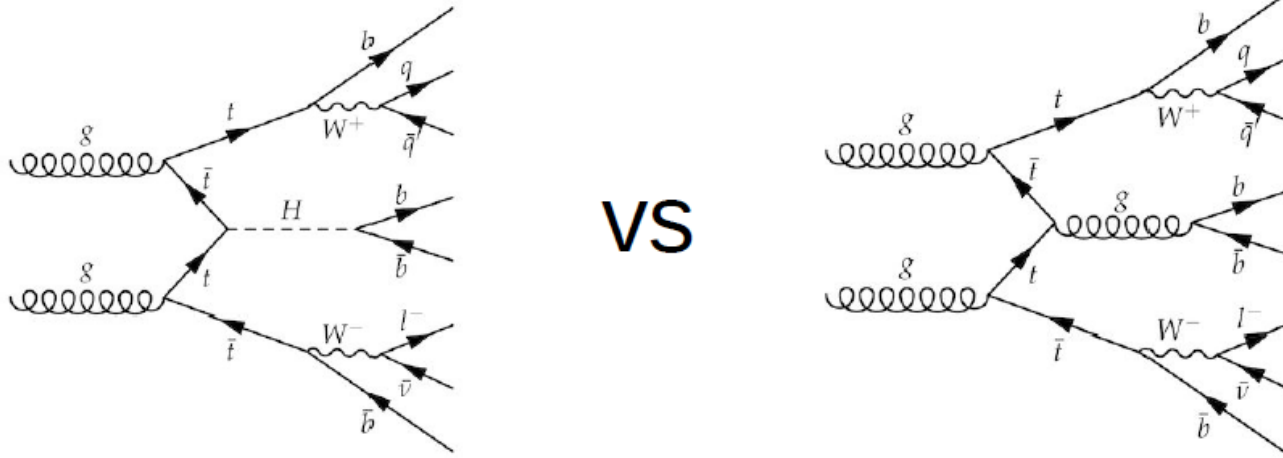
Lorentz boost network: feature extraction

- ◆ Extract generic features from boosted particles.
 - ◆ Single particle features: E, m, pT, η, ϕ .
 - ◆ Pairwise features: such as $\cos(\theta)$ between all pairs.
- ◆ E.g. $\cos(\theta^*)$ in the semi-leptonic decay of the top quark, defined as the angular difference between the direction of the charged lepton in the W rest frame and the direction of the W in the top rest frame.

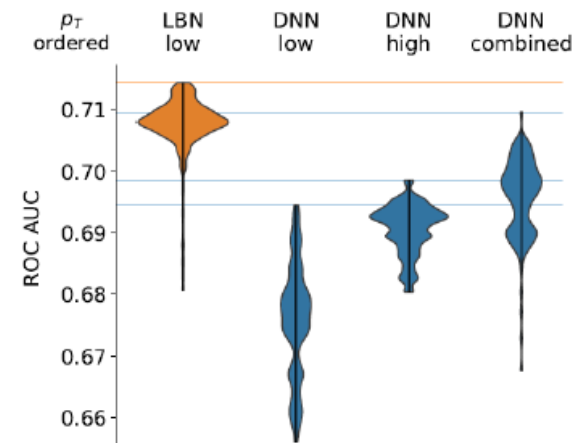


BN for $t\bar{t}(bb)$ vs $t\bar{t}+bb$: performance

- ◆ LBN performance compared to standard DNN with low-, high- and combination of low- and high-level variables.



- ◆ LBN shows improved performance in terms of ROC AUC.



Conclusions

- ◆ Neural networks are widely used in HEP and will become more and more important.
- ◆ A quick overview of the basic structure of the most used NNs in HEP was given.
- ◆ New NNs layers, specifically engineered for HEP, were created.
 - ◆ In this case, high performance comes also with a good interpretability of the trained parameters.
- ◆ NN is a quickly developing field. Exciting time to work on it and to find new applications for HEP.