



Multivariate Data Analysis and Machine Learning in High Energy Physics (III)

Helge Voss (MPI-K, Heidelberg)

Graduierten-Kolleg , Freiburg, 11.5-15.5, 2009



Outline

- Summary of last lecture
- 1-dimensional Likelihood:
 - data preprocessing → decorrelation
- Recap of Fisher's Discriminant
- Neural Networks

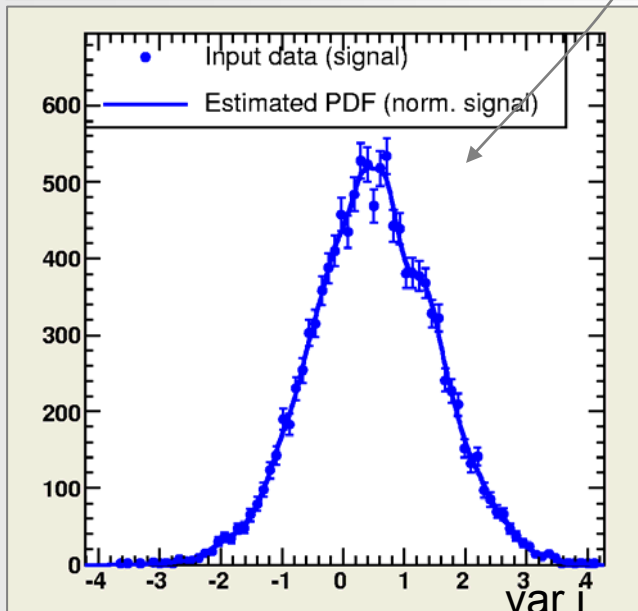
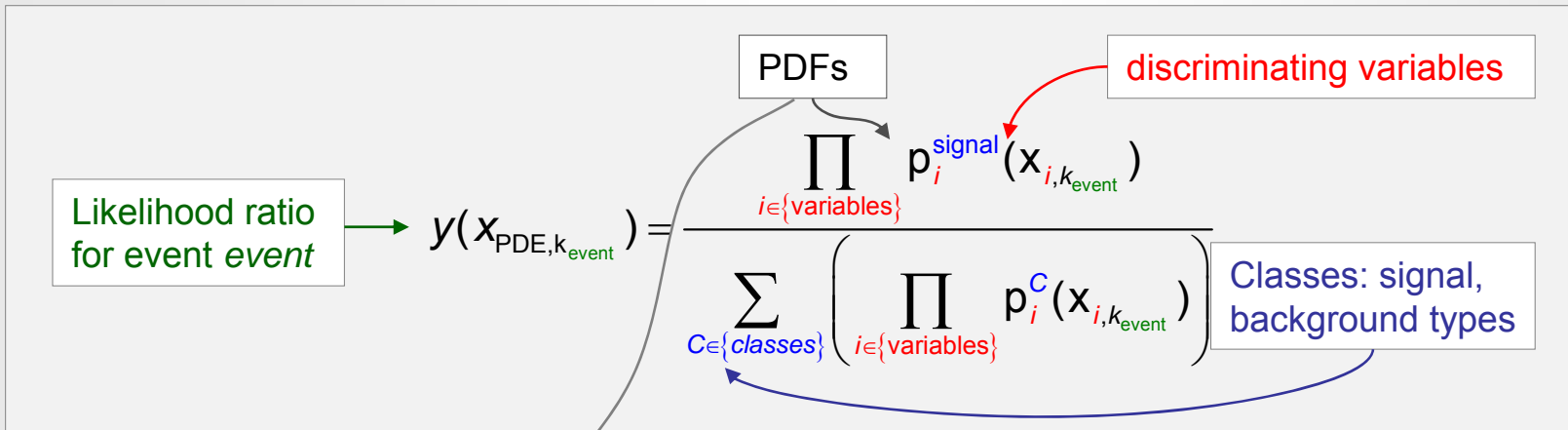
What we've learned yesterday:

- Neyman Pearson: Actually he states the “obvious” if you think about it:
 - if you want to know to which probability an event with measured variables “ x ” is either signal or background, use the ratio between the “true probability” between observing an event with “ x ” from either signal or background.
 - Who is surprised that this is the “best you can do” ?
- Unfortunately we basically never have access to the “true probability density”:
 - hence we try to either estimate the $\text{PDF}_S(x)$ and $\text{PDF}_B(x)$ using Kernel Density Estimator/Multidimensional likelihood (effectively averages of the PDF over regions of the variable space, derived from training events) → curse of dimensionality!!
 - or neglect correlations and use 1-dimensional PDF's → no problems with dimensionality
 - or try a different approach of directly determining hyperplanes in the feature space that separate signal and background events.
- There is no magic, you still need to “tune” parameters (e.g size of the region you want to average over in PDF estimate) in order to get good results.

Naïve Bayesian Classifier

“often called: (projective or 1D) Likelihood”

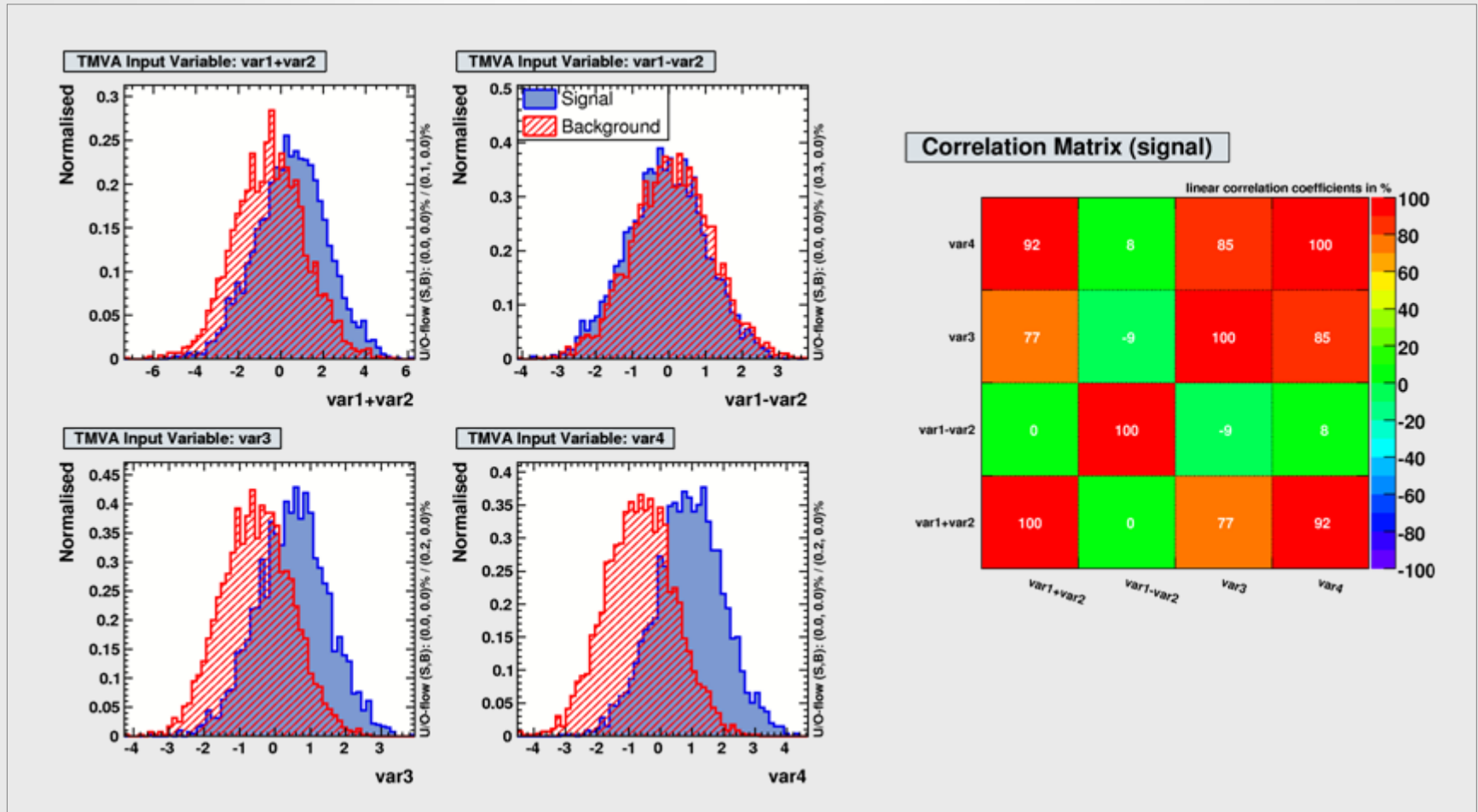
If correlations between variables are weak: $\rightarrow p(\mathbf{x}) \cong \prod_{i=0}^D p_i(\mathbf{x})$



- If the correlations between variables is really negligible, this classifier is “perfect” (simple, robust, performing)
- If not, you seriously loose performance ☹️
- How can we “fix” this ?

What if there are correlations?

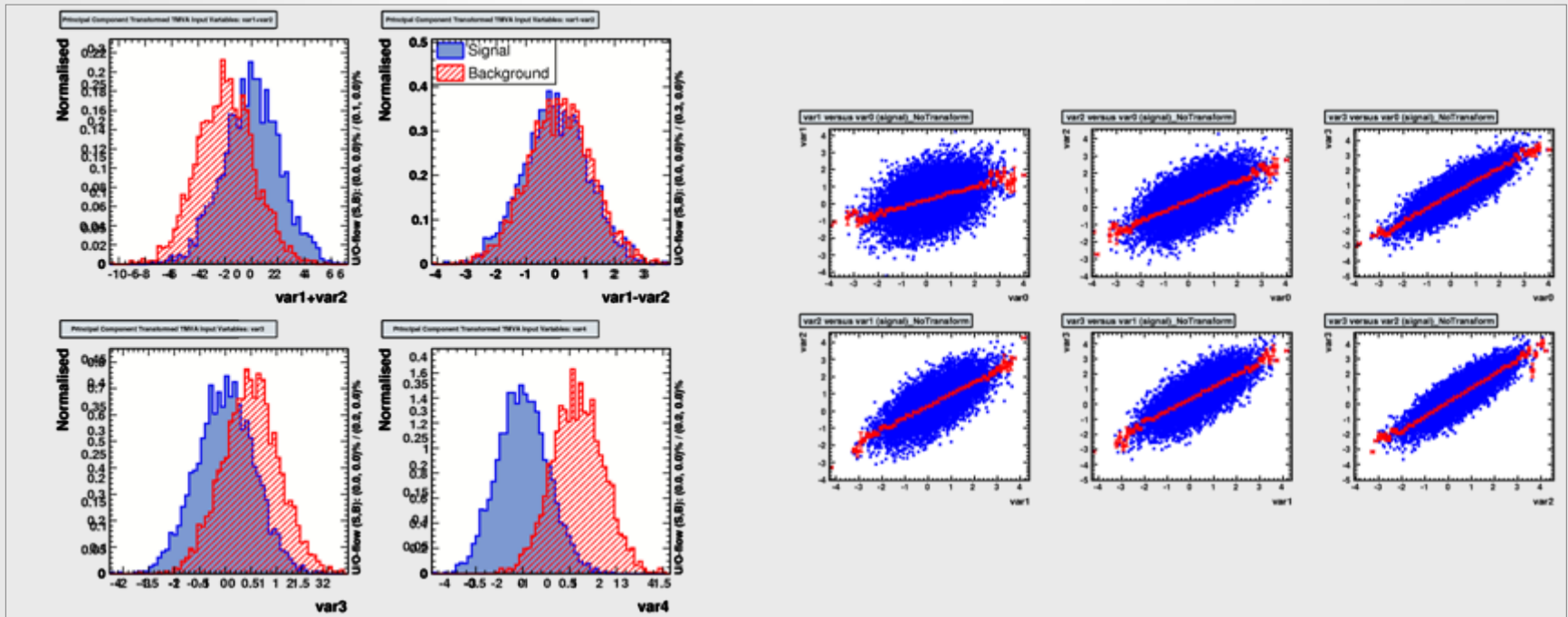
- Typically correlations are present: $C_{ij} = \text{cov}[x_i, x_j] = E[x_i x_j] - E[x_i]E[x_j] \neq 0$ ($i \neq j$)



→ pre-processing: choose set of linear transformed input variables for which $C_{ij} = 0$ ($i \neq j$)

Decorrelation

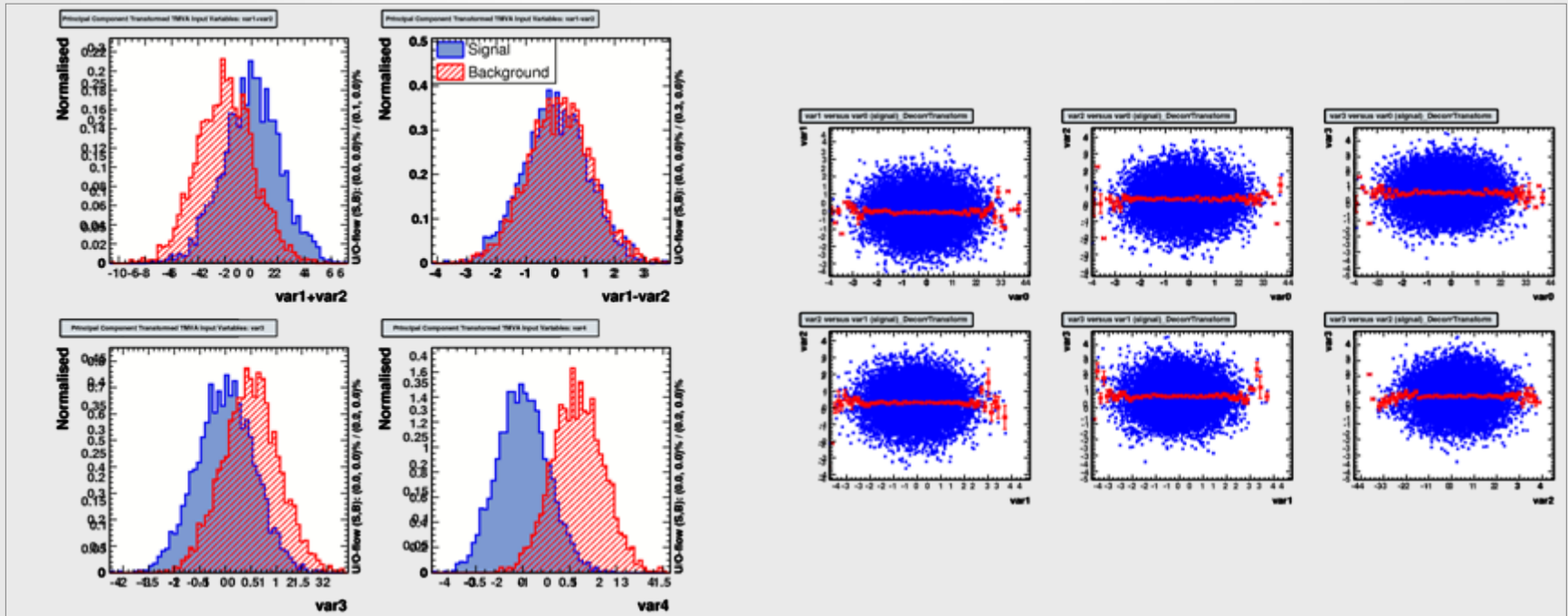
- Determine *square-root* C' of correlation matrix C , i.e., $C = C' C'$
 - compute C' by diagonalising C : $D = S^T C S \Rightarrow C' = S \sqrt{D} S^T$
 - transformation from original (x) in de-correlated variable space (x') by: $x' = C'^{-1} x$



Attention. This is is able to eliminate only linear correlations!!

Decorrelation

- Determine *square-root* C' of correlation matrix C , i.e., $C = C' C'$
 - compute C' by diagonalising C : $D = S^T C S \Rightarrow C' = S \sqrt{D} S^T$
 - transformation from original (x) in de-correlated variable space (x') by: $x' = C'^{-1} x$



Attention. This is is able to eliminate only linear correlations!!

Decorrelation: Principal Component Analysis

- PCA is typically used to:
 - reduce dimensionality of a problem
 - find the most dominant features in your distribution by transforming
- The eigenvectors of the covariance matrix with the largest eigenvalues correspond to the dimensions that have the strongest correlation in the data set. Along these axis the variance is largest
 - sort the eigenvectors according to their eigenvalues
- Dataset is transformed in variable space along these eigenvectors
 - Along the “first” dimension the data show the largest “features”, the smallest features are found in the “last” dimension.

$$x_k^{\text{PC}}(i_{\text{event}}) = \sum_{v \in \{\text{variables}\}} [x_v(i_{\text{event}}) - \bar{x}_v] \cdot v_v^{(k)}, \quad \forall k \in \{\text{variables}\}$$

Principle Component (PC) of variable k

sample means

eigenvector


- Matrix of eigenvectors V obey the relation: $C \cdot V = D \cdot V \rightarrow$ **PCA eliminates correlations!**
 - correlation matrix
 - diagonalised square root of C

How to Apply the Pre-Processing Transformation?

In general: the decorrelation for signal and background variables is different
however: for a “test event” we don’t know beforehand if it is signal or background.

? What do we do?

→ for likelihood ratio, decorrelate signal and background independently

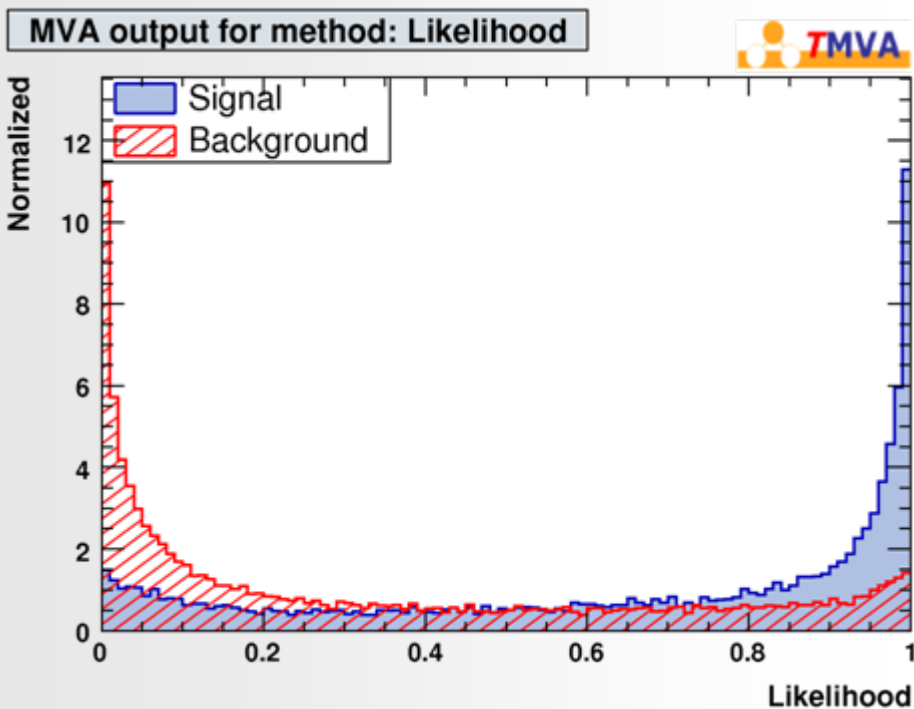
$$y_L(i_{\text{event}}) = \frac{\prod_{k \in \{\text{variables}\}} p_k^S(x_k(i_{\text{event}}))}{\prod_{k \in \{\text{variables}\}} p_k^S(x_k(i_{\text{event}})) + \prod_{k \in \{\text{variables}\}} p_k^B(x_k(i_{\text{event}}))}$$

$$y_L^{\text{trans}}(i_{\text{event}}) = \frac{\prod_{k \in \{\text{variables}\}} p_k^S(\hat{T}^S x_k(i_{\text{event}}))}{\prod_{k \in \{\text{variables}\}} p_k^S(\hat{T}^S x_k(i_{\text{event}})) + \prod_{k \in \{\text{variables}\}} p_k^B(\hat{T}^B x_k(i_{\text{event}}))}$$

→ for other estimators, one need to decide on one of the two...

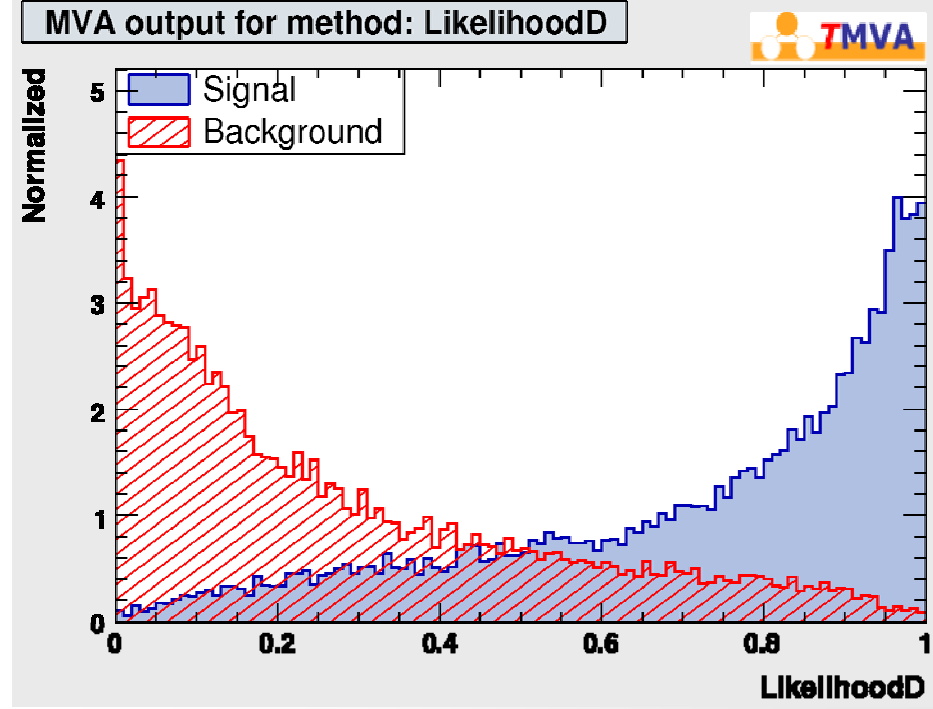
Decorrelation at Work

- Example: linear correlated Gaussians \rightarrow decorrelation works to 100%
- \rightarrow 1-D Likelihood on decorrelated sample give best possible performance
- \rightarrow compare also the effect on the MVA-output variable!

correlated variables:



after decorrelation

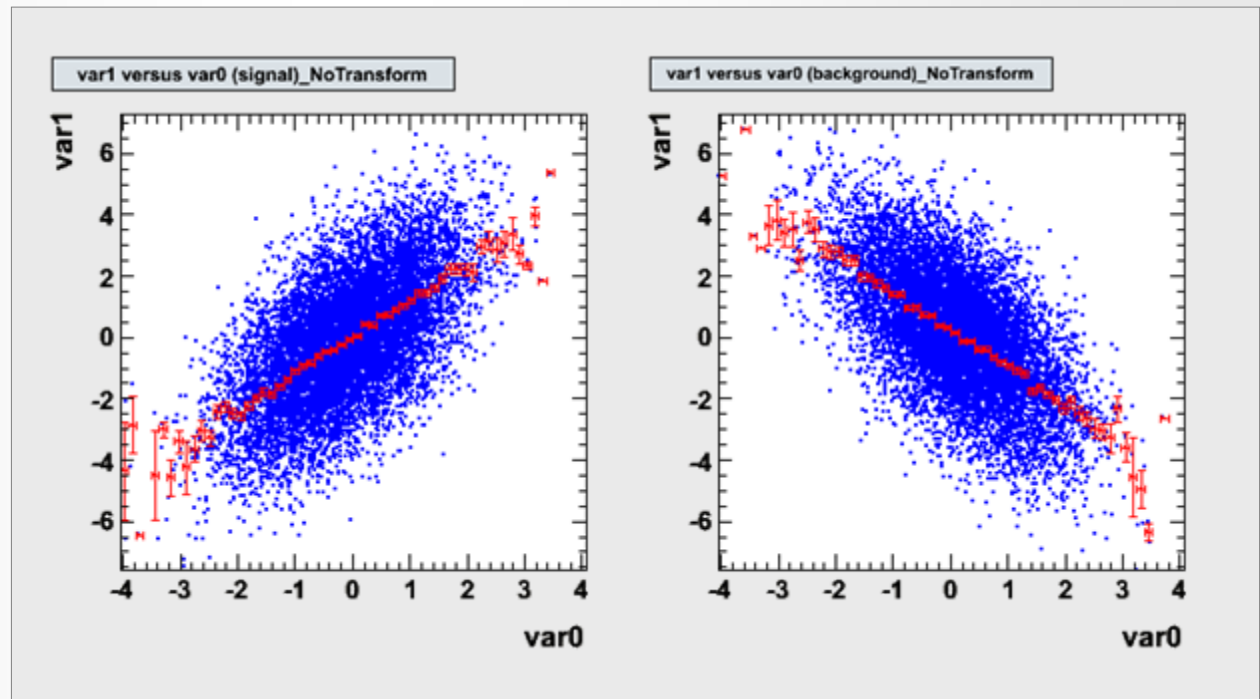


(note the different scale on the y-axis... sorry)

Limitations of the Decorrelation

- in cases with non-Gaussian distributions and/or nonlinear correlations, the decorrelation needs to be treated with care
- How does linear decorrelation affect cases where correlations between signal and background differ?

Original correlations



Signal

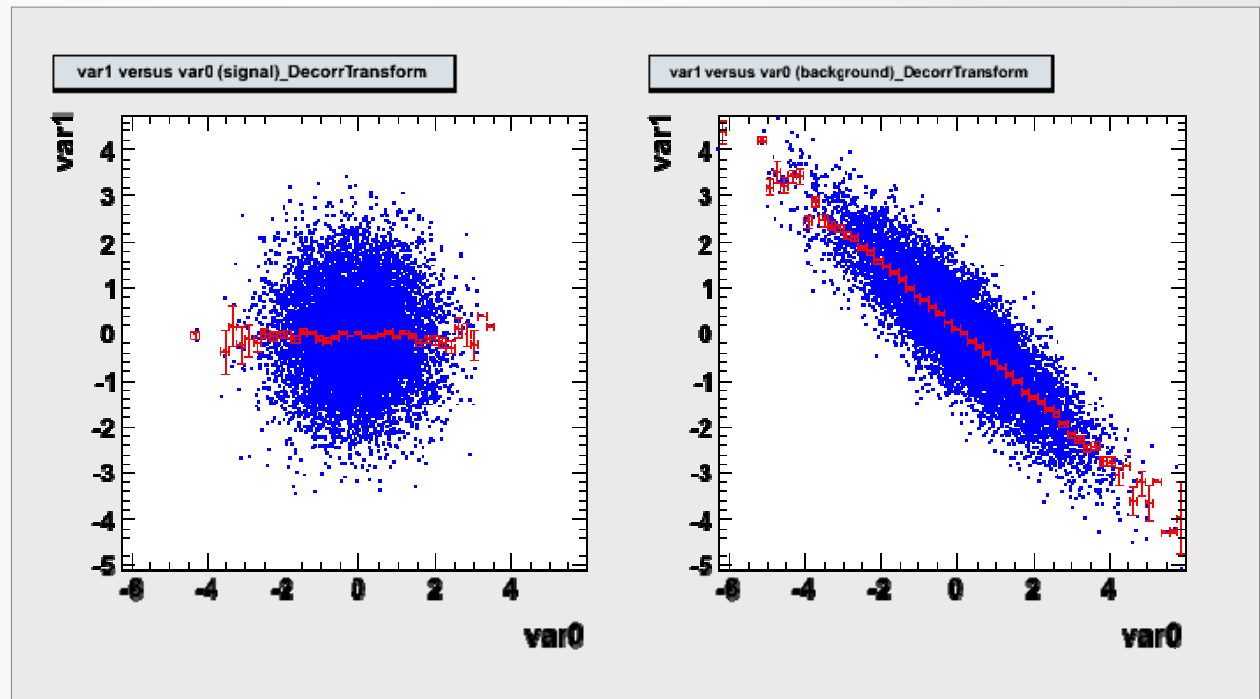
Background

Limitations of the Decorrelation

- in cases with non-Gaussian distributions and/or nonlinear correlations, the decorrelation needs to be treated with care
- How does linear decorrelation affect cases where correlations between signal and background differ?

Original correlations

SQRT decorrelation



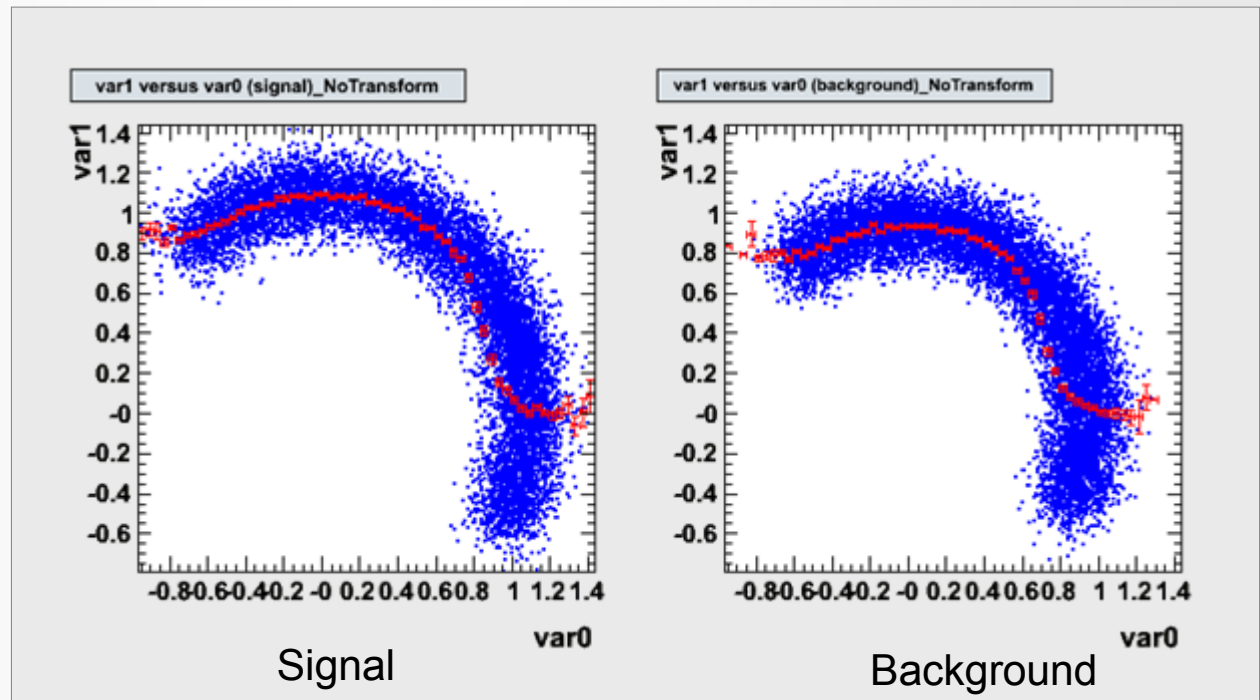
Signal

Background

Limitations of the Decorrelation

- in cases with non-Gaussian distributions and/or nonlinear correlations, the decorrelation needs to be treated with care
- How does linear decorrelation affect strongly nonlinear cases ?

Original correlations

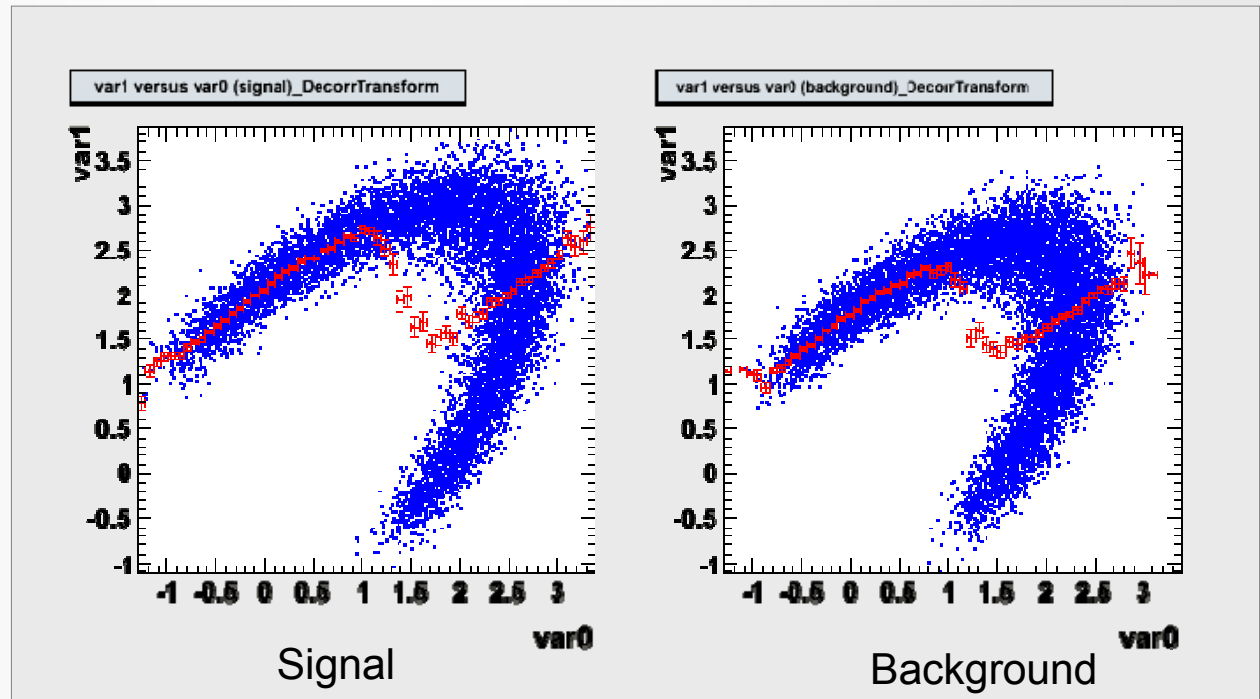


Limitations of the Decorrelation

- in cases with non-Gaussian distributions and/or nonlinear correlations, the decorrelation needs to be treated with care
- How does linear decorrelation affect strongly nonlinear cases ?

Original correlations

SQRT decorrelation

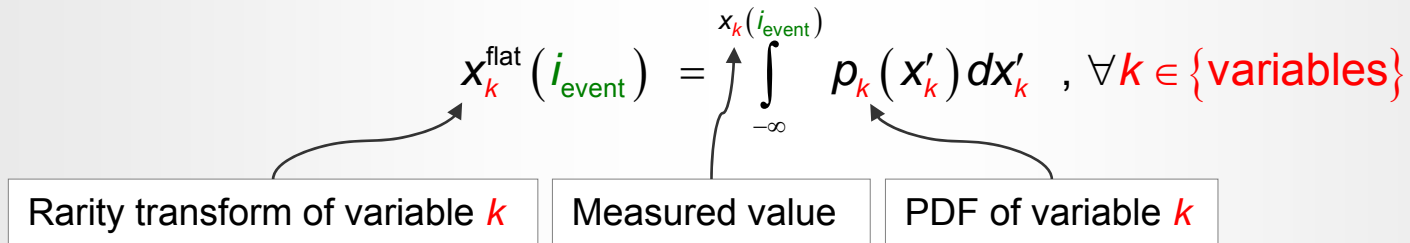


→ Watch out before you used decorrelation “blindly”!!

“Gaussian-isation”

- Improve decorrelation by pre-Gaussianisation of variables

- ➔ First: transformation to achieve uniform (flat) distribution:

$$x_k^{\text{flat}}(i_{\text{event}}) = \int_{-\infty}^{x_k(i_{\text{event}})} p_k(x'_k) dx'_k, \quad \forall k \in \{\text{variables}\}$$


Rarity transform of variable k Measured value PDF of variable k

The integral can be solved in an unbinned way by event counting, or by creating non-parametric PDFs (see later for likelihood section)

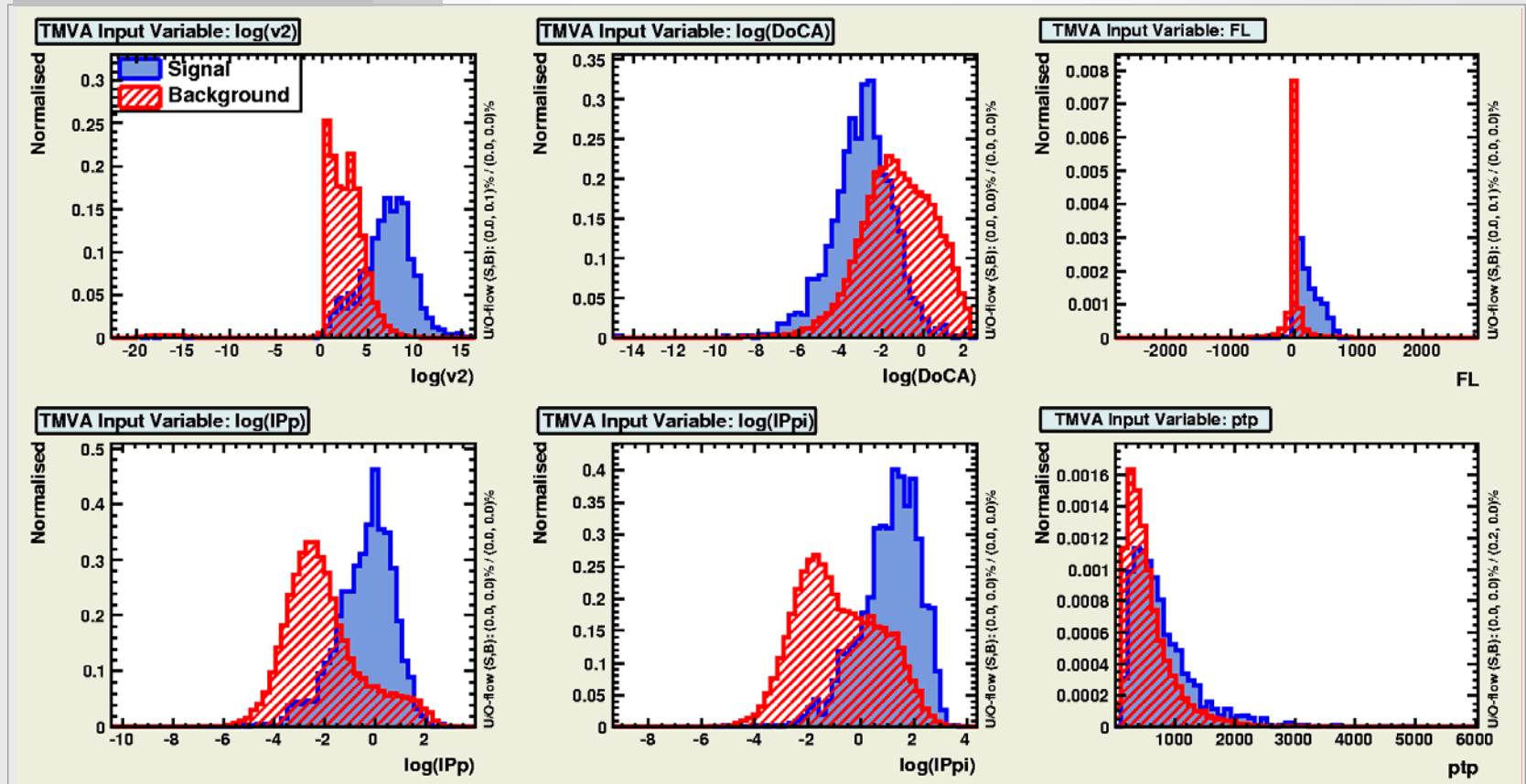
- ➔ Second: make Gaussian via inverse error function: $\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$

$$x_k^{\text{Gauss}}(i_{\text{event}}) = \sqrt{2} \cdot \text{erf}^{-1}(2x_k^{\text{flat}}(i_{\text{event}}) - 1), \quad \forall k \in \{\text{variables}\}$$

- ➔ Third: decorrelate (and “iterate” this procedure)

“Gaussian-isation”

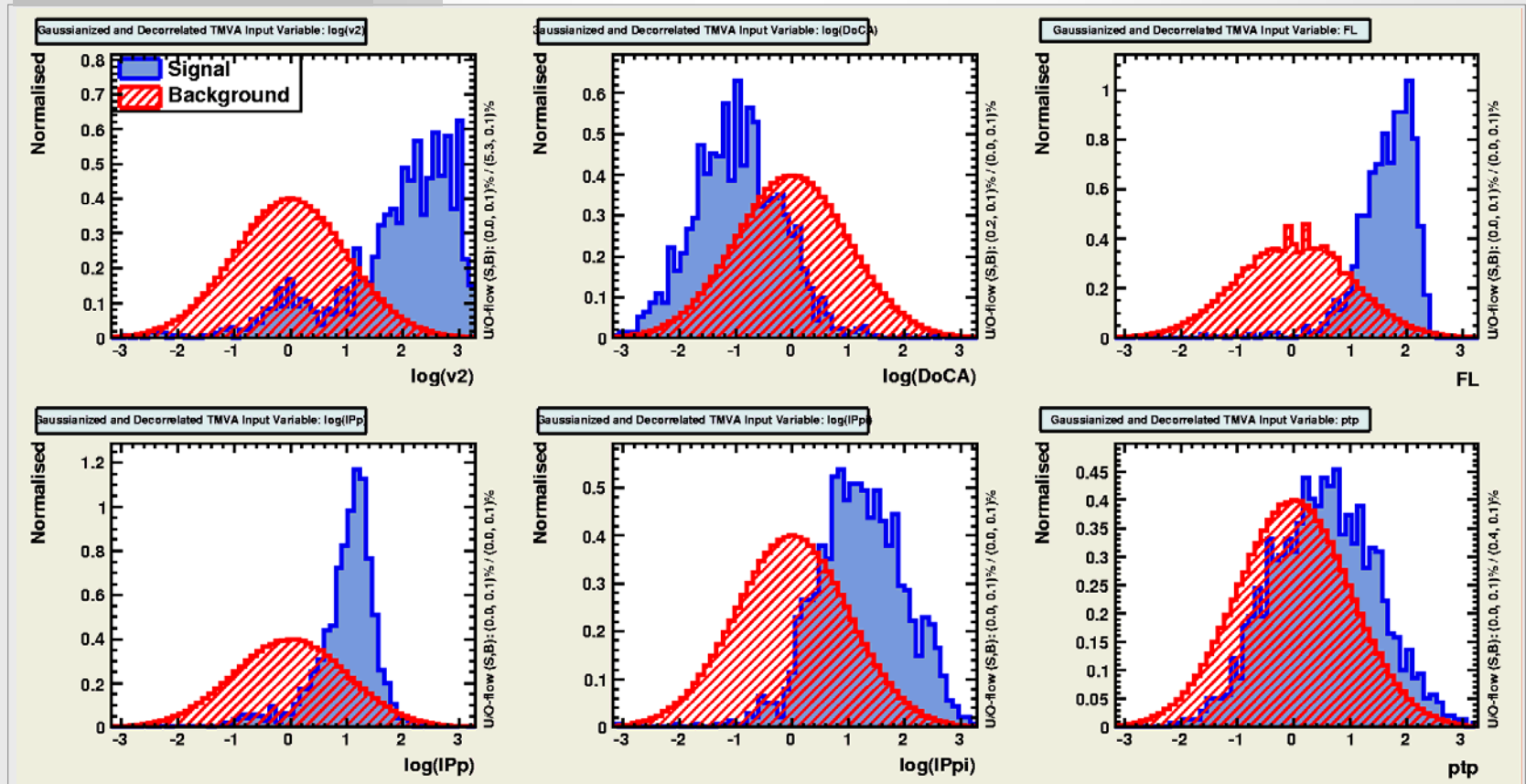
Background - Gaussianised



We cannot simultaneously Gaussianise both signal and background ?!

“Gaussian-isation”

Background - Gaussianised



We cannot simultaneously Gaussianise both signal and background ?!

Linear Discriminant

If non parametric methods like ‘k-Nearest-Neighbour’ or ‘Kernel Density Estimators’ suffer from

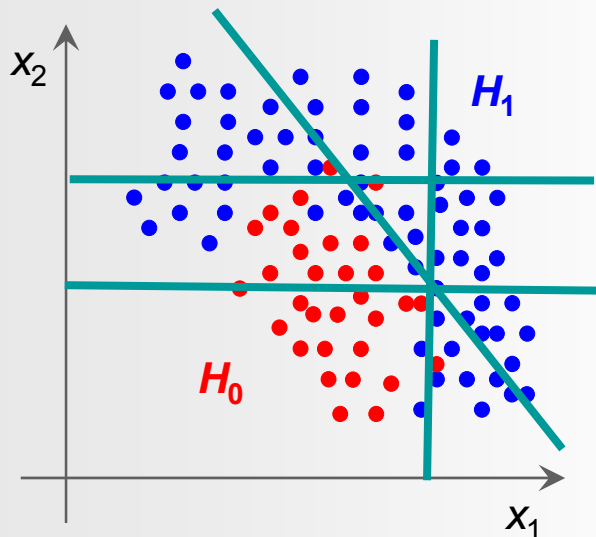
- lack of training data → “curse of dimensionality”
- slow response time → need to evaluate the whole training data for every test event

→ use of parametric models $y(x)$ to fit to the training data
$$y(x = \{x_1, \dots, x_D\}) \approx \sum_{i=0}^M w_i h_i(x)$$

Linear Discriminant:

i.e. any linear function of the input variables:
giving rise to linear decision boundaries

$$y(x = \{x_1, \dots, x_D\}) \approx w_0 + \sum_{i=1}^D w_i x_i$$

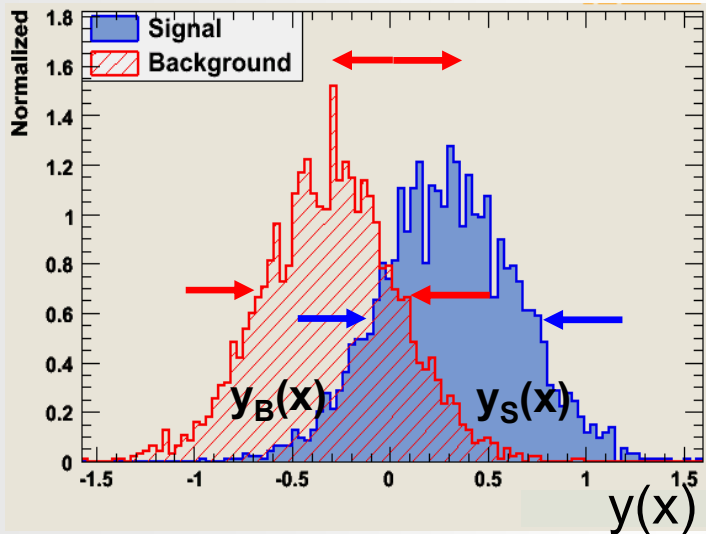


How do we determine the “weights” w that do “best”??

Watch out: these lines are NOT the linear function $y(x)$, but rather the decision boundary given by $y(x)=\text{const}$.

Fisher's Linear Discriminant

$$y(x = \{x_1, \dots, x_D\}) = y(x, w) = w_0 + \sum_{i=1}^D w_i x_i$$



How do we determine the “weights” w that do “best”??

- Maximise the “separation” between the two classes S and B
 - minimise overlap of the distribution $y_S(x)$ and $y_B(x)$
 - maximise the distance between the two mean values of the classes
 - minimise the variance within each class

→ maximise
$$J(\vec{w}) = \frac{(E(y_B) - E(y_S))^2}{\sigma_{y_B}^2 + \sigma_{y_S}^2} = \frac{\vec{w}^T B \vec{w}}{\vec{w}^T W \vec{w}} = \frac{\text{"in between" variance}}{\text{"within" variance}}$$

$$\vec{\nabla}_{\vec{w}} J(\vec{w}) = 0 \Rightarrow \vec{w} \propto W^{-1} (\langle \vec{x} \rangle_S - \langle \vec{x} \rangle_B) ; \quad \text{the Fisher coefficients}$$

note: these quantities can be calculated from the training data

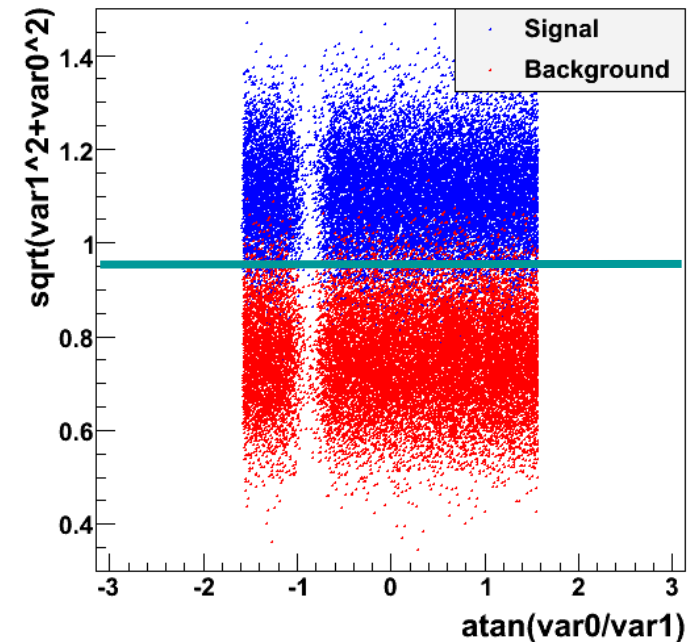
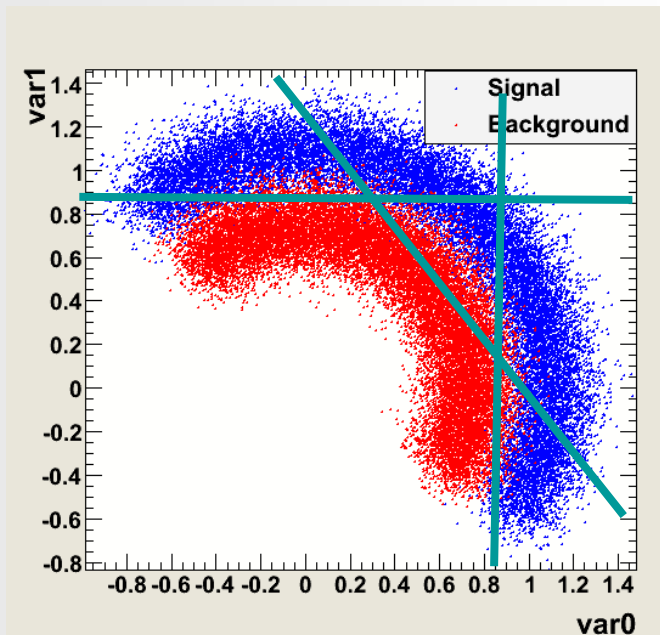
Linear Discriminant and non linear correlations

assume the following non-linear correlated data:

- the Linear discriminant obviously doesn't do a very good job here:
- Of course, these can easily be de-correlated:
 - here: linear discriminator works perfectly on de-correlated data

$$\text{var } 0^l = \sqrt{\text{var } 0^2 + \text{var } 1^2}$$

$$\text{var } 1^l = \text{atan}\left(\frac{\text{var } 0}{\text{var } 1}\right)$$



Linear Discriminant with Quadratic input:

- A simple extension of the linear decision boundary to “quadratic” decision boundary:

while: ■ var0
■ var1

→ linear decision boundaries in var0,var1

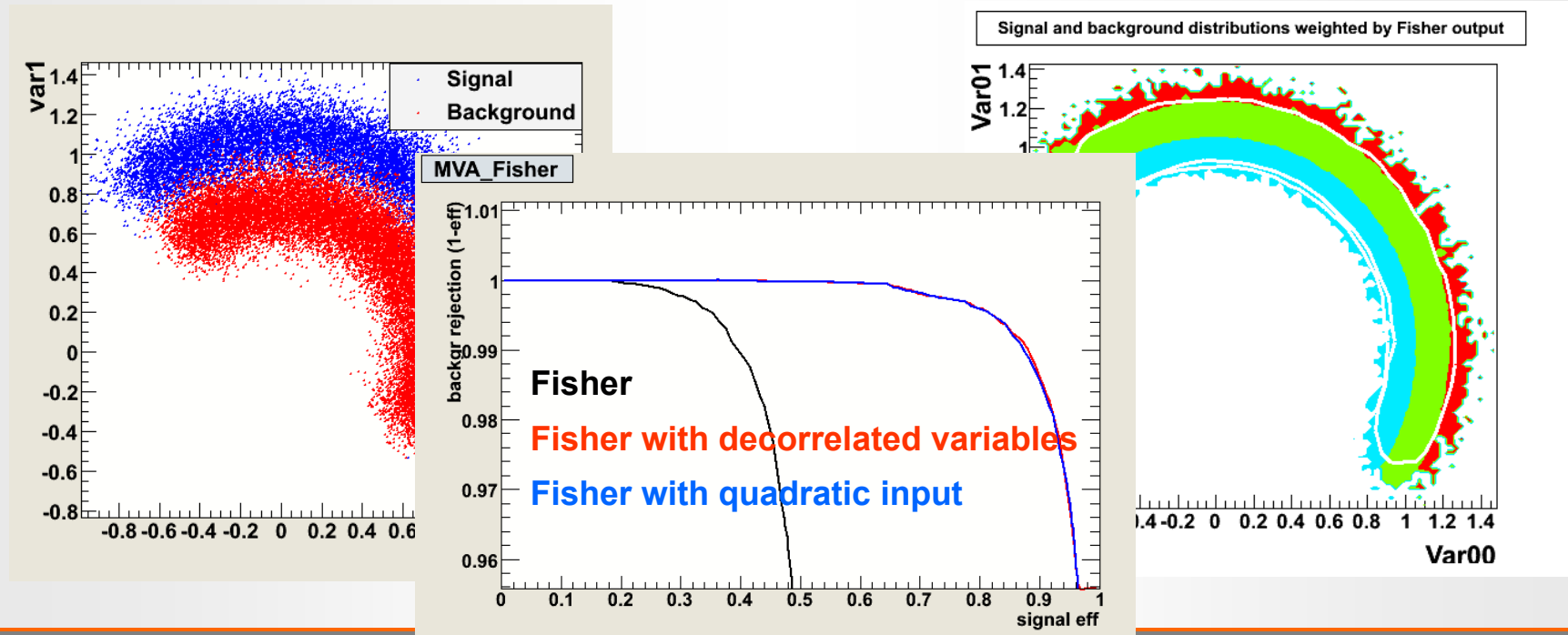
■ var0 * var0

■ var1 * var1

■ var0 * var1

→ quadratic decision boundaries in var0,var1

Performance of Fisher Discriminant:
with quadratic input:



Training Classifiers and Loss-Function

- All classifiers (KNN, Likelihood, Fisher) could be calculated and didn't require parameter fitting
- Other classifiers provide a set of "basis" functions (or model) that need to be optimally adjusted to find the appropriate separating hyperplane (surface)
- Let $x \in \mathbb{R}^n$ be a random variable (i.e. our observables) and $y(x)$
 - y a real valued output variable with joint distribution $\Pr(x, y) \rightarrow$ regression
 - y value $\rightarrow 1$ for **signal**, $y \rightarrow 0$ for **background** \rightarrow classification
- we are looking for a function $y(x)$ that predicts y given certain input variables:
 $y(x): \mathbb{R}^n \rightarrow \mathbb{R}$:
- Loss function: $L(y, y(x))$ penalizes errors made in the prediction:
 - $EPE(y(x)) = E(y - y(x))^2$ squared error loss, typical "loss function" used in regression
 - by conditioning on "x" we can write:
 - $EPE(y(x)) = E(|y - y(x)|)$ misclassification error, typical "loss function" for classification problems

Neural Networks

naturally, if we want to go the “arbitrary” non-linear decision boundaries, $y(x)$ needs to be constructed in “any” non-linear fashion

$$y(\vec{x}) = \sum_i^M (w_i h_i(\vec{x}))$$

- Think of $h_i(x)$ as a set of “basis” functions
- If $h(x)$ is sufficiently general (i.e. non linear), a linear combination of “enough” basis function should allow to describe any possible discriminating function $y(x)$

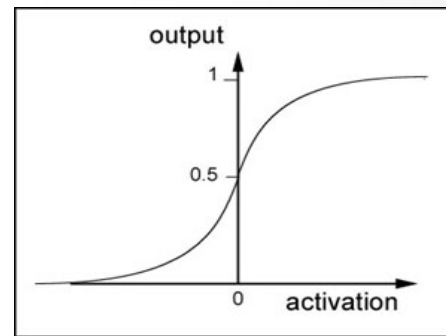
K.Weierstrass Theorem: proves just that previous statement.

Imagine you chose $h_i(x)$ to be such that:

$$y(x) = \sum_i^M w_{0i} A \left(w_{i0} + \sum_{j=1}^D w_{ij} \cdot x_j \right)$$

$y(x) =$
a linear combination of
non linear functions of
linear combination of
the input data

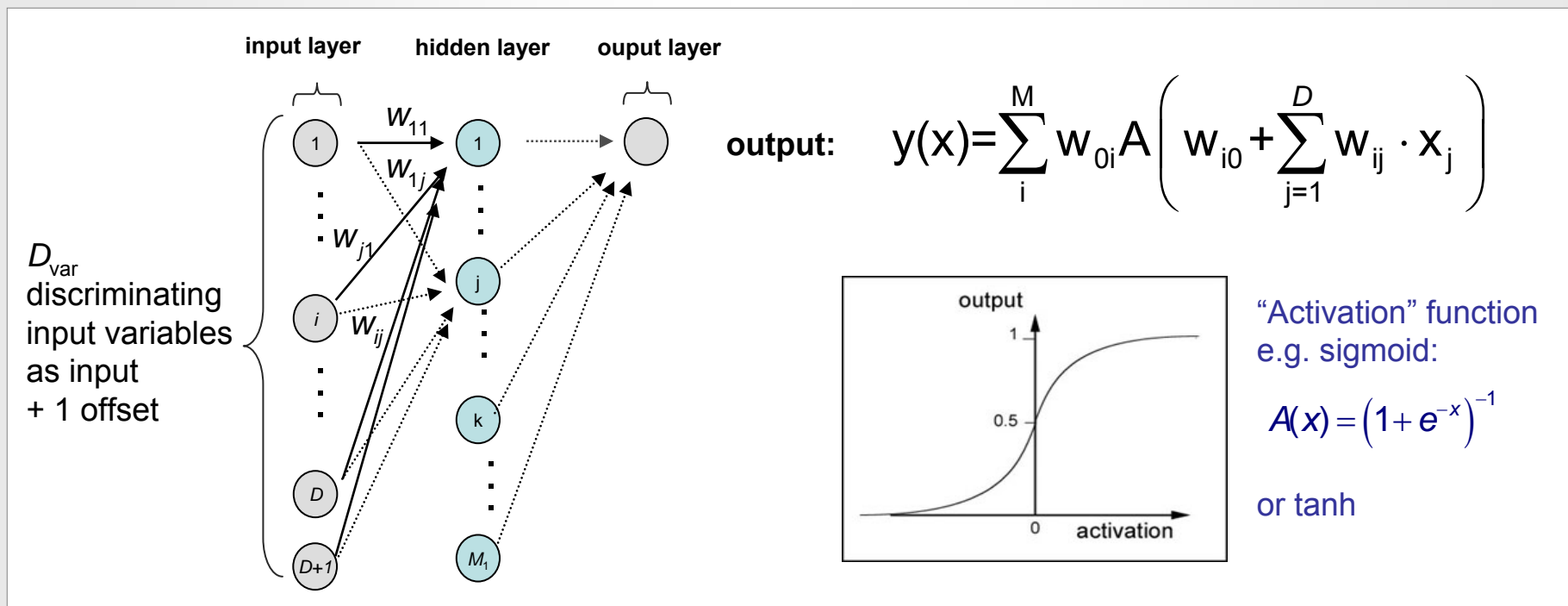
$$A(x) = \frac{1}{1+e^{-x}} : \text{the sigmoid function}$$



Now we “only” need to find the appropriate “weights” w

Neural Networks: Multilayer Perceptron MLP

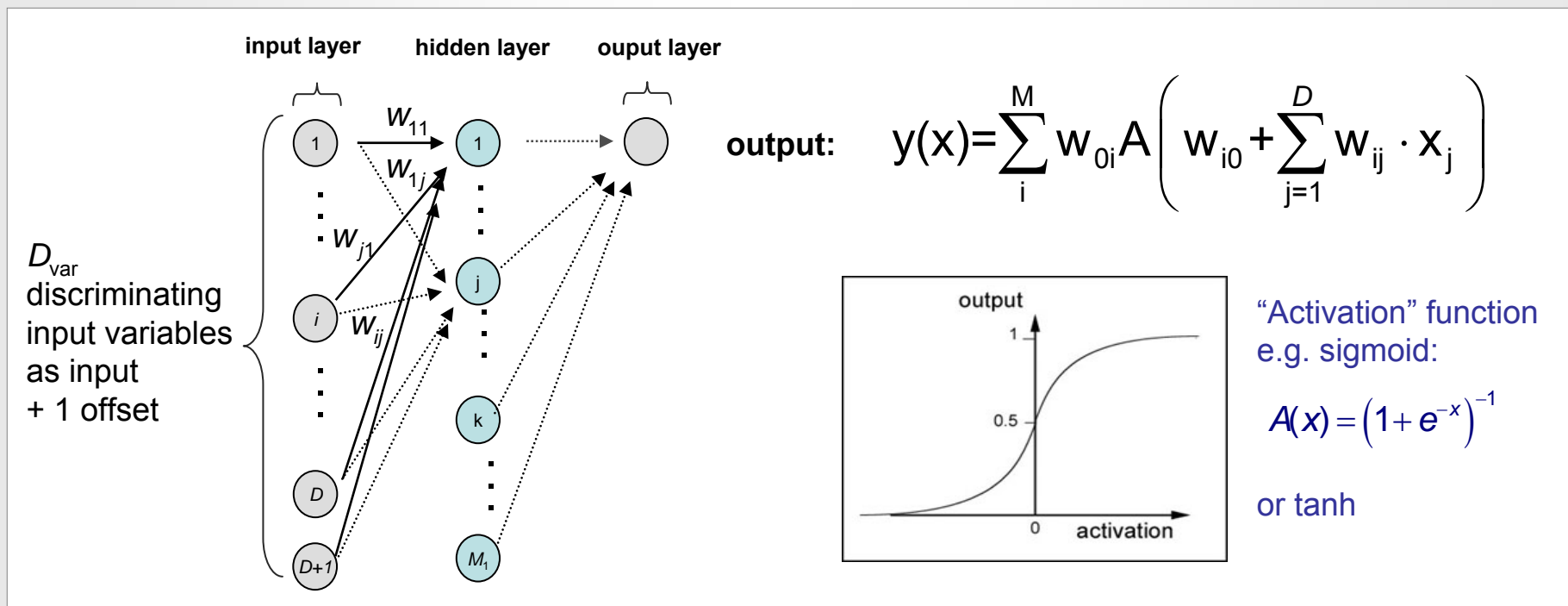
But before talking about the weights, let's try to "interpret" the formula as a Neural Network:



- Nodes in hidden layer represent the “activation functions” whose arguments are linear combinations of input variables → non-linear response to the input
- The output is a linear combination of the output of the activation functions at the internal nodes
- Input to the layers from preceding nodes only → feed forward network (no backward loops)
- It is straightforward to extend this to “several” input layers

Neural Networks: Multilayer Perceptron MLP

But before talking about the weights, let's try to "interpret" the formula as a Neural Network:



nodes → neurons
 links(weights) → synapses



Neural network: try to simulate reactions of a brain to certain stimulus (input data)

Neural Network Training

idea: using the “training events” adjust the weights such, that

- $y(x) \rightarrow 0$ for background events
- $y(x) \rightarrow 1$ for signal events

how do we adjust ?

- minimize Loss function:

$$L(w) = \sum_i^{\text{events}} \underbrace{(y(x_i))}_{\text{predicted event type}} - \underbrace{y(C)}_{\text{true event type}})^2$$

where

i.e. use usual “sum of squares”

$$y(C) = \begin{cases} 1 & \text{for } C = \text{signal} \\ 0 & \text{for } C = \text{backgr.} \end{cases}$$

- $y(x)$ is a very “wiggly” function with many local minima. A global overall fit in the many parameters is possible but not the most efficient method to train neural networks

- back propagation (learn from experience, gradually adjust your perception to match reality)
- online learning (learn event by event -- not only at the end of your life from all experience)

Neural Network Training

back-propagation and online learning

- start with random weights
- adjust weights in each step a bit in the direction of the steepest descend of the “Loss”-function L

$$\mathbf{w}^{n+1} = \mathbf{w}^n + \eta \cdot \overrightarrow{\nabla}_{\mathbf{w}} L(\mathbf{w}) \quad (\mathbf{w}) \quad \eta = \text{learning rate} \quad L(\mathbf{w}) = (y(\mathbf{x}_i) - y(\mathbf{C}))^2$$

- for weights connected to output nodes

$$\frac{\partial L}{\partial \mathbf{w}_{oi}} = (y(\mathbf{x}) - y(\mathbf{C})) A \left(\mathbf{w}_{i0} + \sum_{j=1}^D \mathbf{w}_{ij} \cdot \mathbf{x}_j \right)$$

$$y(\mathbf{x}) = \sum_i^M \mathbf{w}_{oi} A \left(\mathbf{w}_{i0} + \sum_{j=1}^D \mathbf{w}_{ij} \cdot \mathbf{x}_j \right)$$

- for weights connected to output nodes
... a bit more complicated formula

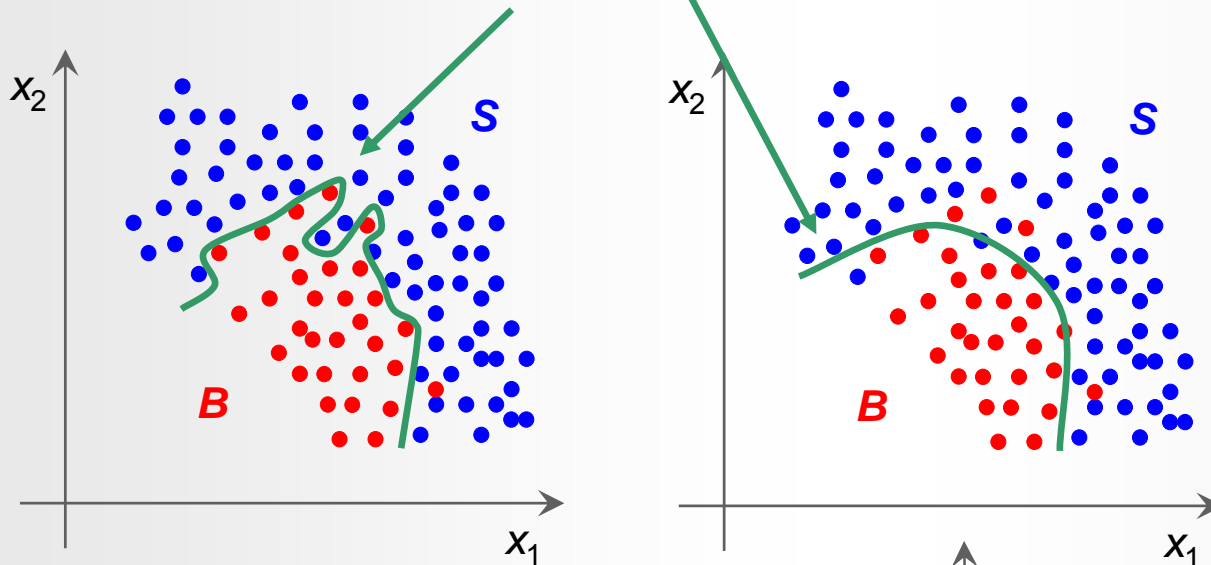
- note: all these gradients are easily calculated from the training event

- training is repeated n-times over the whole training data sample. **how often ??**

for online learning, the training events should be mixed randomly, otherwise you first steer in a wrong direction from which it is afterwards hard to get out again!

Overtraining

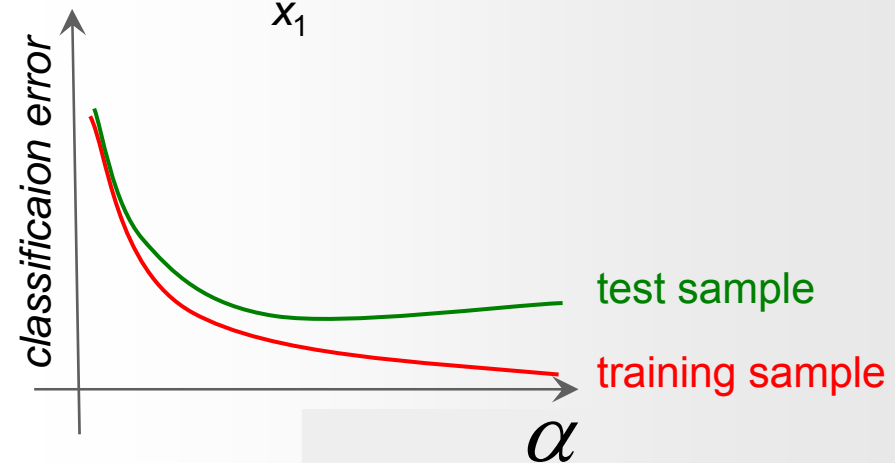
- training is repeated n-times over the whole training data sample. **how often ??**
- it seems intuitive that this boundary will give better results in another statistically independent data set than that one



→ stop training before you learn statistical fluctuations in the data

→ verify on independent “test” sample

→ possible overtraining is concern for every “tunable parameter” α of classifiers:
Smoothing parameter, n-nodes...



Cross Validation

- many (all) classifiers have tuning parameters “ α ” that need to be controlled against overtraining
 - number of training cycles, number of nodes (neural net)
 - smoothing parameter h (kernel density estimator)
 -
- the more free parameter a classifier has to adjust internally \rightarrow more prone to overtraining
- more training data \rightarrow better training results
- division of data set into “training” and “test” sample reduces the “training” data

Cross Validation: divide the data sample into say 5 sub-sets

Train

Train

Train

Train

Test

- train 5 classifiers: $y_i(x, \alpha) : i=1, \dots, 5$,
- classifier $y_i(x, \alpha)$ is trained without the i -th sub sample
- calculate the test error:
$$CV(\alpha) = \frac{1}{N_{\text{events}}} \sum_k^{\text{events}} L(y_i(x_k, \alpha)) \quad L : \text{loss function}$$
- choose tuning parameter α for which $CV(\alpha)$ is minimum and train the final classifier using all data

What is the best Network Architecture?

- Theoretically a single hidden layer is enough for any problem, provided one allows for sufficient number of nodes. (K.Weierstrass theorem)
- “Relatively little is known concerning advantages and disadvantages of using a single hidden layer with many nodes over many hidden layers with fewer nodes. The mathematics and approximation theory of the MLP model with more than one hidden layer is not very well understood

....”nonetheless there seems to be reason to conjecture that the two hidden layer model may be significantly more promising than the single hidden layer model”

A.Pinkus, “Approximation theory of the MLP model with neural networks”,
Acta Numerica (1999),pp.143-195

(Glen Cowan)

- Typically in high-energy physics, non-linearities are reasonably simple,
 - in many cases 1 layer with a larger number of nodes should be enough
 - but it might well be worth trying more layers (and less nodes in each layer)