

# Statistical Methods for Particle Physics

## Lecture 2: hypothesis tests I; multivariate methods

<https://indico.weizmann.ac.il//conferenceDisplay.py?confId=52>



Statistical Inference for Astro  
and Particle Physics Workshop  
Weizmann Institute, Rehovot  
March 8-12, 2015



Glen Cowan  
Physics Department  
Royal Holloway, University of London  
[g.cowan@rhul.ac.uk](mailto:g.cowan@rhul.ac.uk)  
[www.pp.rhul.ac.uk/~cowan](http://www.pp.rhul.ac.uk/~cowan)

# Outline for Monday – Thursday

(GC = Glen Cowan, KC = Kyle Cranmer)

## Monday 9 March

GC: probability, random variables and related quantities  
KC: parameter estimation, bias, variance, max likelihood

## Tuesday 10 March

KC: building statistical models, nuisance parameters  
GC: hypothesis tests I,  $p$ -values, multivariate methods

## Wednesday 11 March

KC: hypothesis tests 2, composite hyp., Wilks', Wald's thm.  
GC: asymptotics 1, Asimov data set, sensitivity

## Thursday 12 March:

KC: confidence intervals, asymptotics 2  
GC: unfolding

# Hypotheses

A hypothesis  $H$  specifies the probability for the data, i.e., the outcome of the observation, here symbolically:  $x$ .

$x$  could be uni-/multivariate, continuous or discrete.

E.g. write  $x \sim f(x|H)$ .

$x$  could represent e.g. observation of a single particle, a single event, or an entire “experiment”.

Possible values of  $x$  form the sample space  $S$  (or “data space”).

Simple (or “point”) hypothesis:  $f(x|H)$  completely specified.

Composite hypothesis:  $H$  contains unspecified parameter(s).

The probability for  $x$  given  $H$  is also called the likelihood of the hypothesis, written  $L(x|H)$ .

# Definition of a test

Goal is to make some statement based on the observed data  $x$  as to the validity of the possible hypotheses.

Consider e.g. a simple hypothesis  $H_0$  and alternative  $H_1$ .

A **test** of  $H_0$  is defined by specifying a **critical region**  $W$  of the data space such that there is no more than some (small) probability  $\alpha$ , assuming  $H_0$  is correct, to observe the data there, i.e.,

$$P(x \in W | H_0) \leq \alpha$$

If  $x$  is observed in the critical region, reject  $H_0$ .

$\alpha$  is called the **size** or **significance level** of the test.

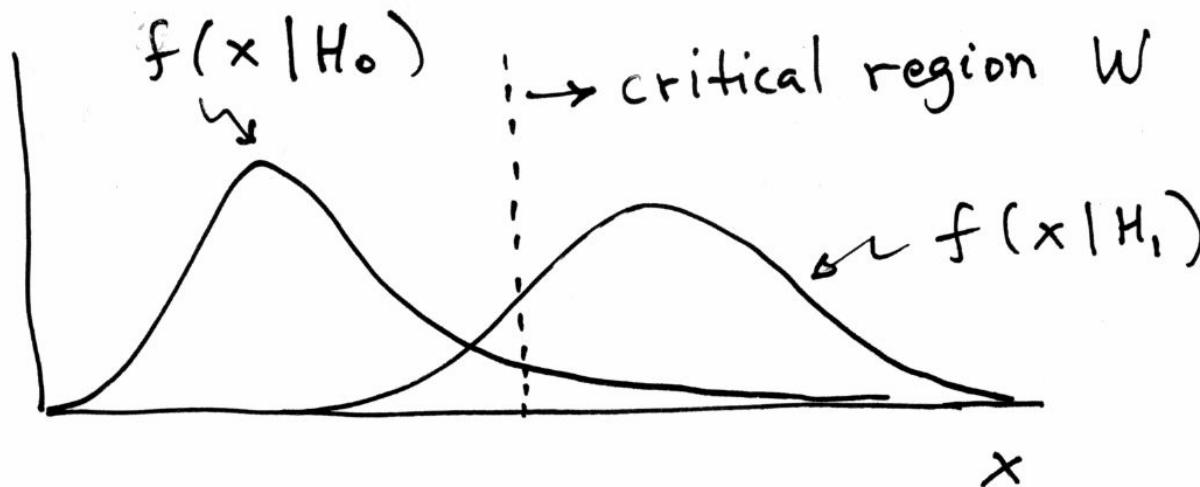
Critical region also called “rejection” region; complement is acceptance region.

## Definition of a test (2)

But in general there are an infinite number of possible critical regions that give the same significance level  $\alpha$ .

So the choice of the critical region for a test of  $H_0$  needs to take into account the alternative hypothesis  $H_1$ .

Roughly speaking, place the critical region where there is a low probability to be found if  $H_0$  is true, but high if  $H_1$  is true:



# Rejecting a hypothesis

Note that rejecting  $H_0$  is not necessarily equivalent to the statement that we believe it is false and  $H_1$  true. In frequentist statistics only associate probability with outcomes of repeatable observations (the data).

In Bayesian statistics, probability of the hypothesis (degree of belief) would be found using Bayes' theorem:

$$P(H|x) = \frac{P(x|H)\pi(H)}{\int P(x|H)\pi(H) dH}$$

which depends on the prior probability  $\pi(H)$ .

What makes a frequentist test useful is that we can compute the probability to accept/reject a hypothesis assuming that it is true, or assuming some alternative is true.

## Type-I, Type-II errors

Rejecting the hypothesis  $H_0$  when it is true is a Type-I error.

The maximum probability for this is the size of the test:

$$P(x \in W | H_0) \leq \alpha$$

But we might also accept  $H_0$  when it is false, and an alternative  $H_1$  is true.

This is called a Type-II error, and occurs with probability

$$P(x \in S - W | H_1) = \beta$$

One minus this is called the power of the test with respect to the alternative  $H_1$ :

$$\text{Power} = 1 - \beta$$

# Choosing a critical region

To construct a test of a hypothesis  $H_0$ , we can ask what are the relevant alternatives for which one would like to have a high power.

Maximize power wrt  $H_1$  = maximize probability to  
reject  $H_0$  if  $H_1$  is true.

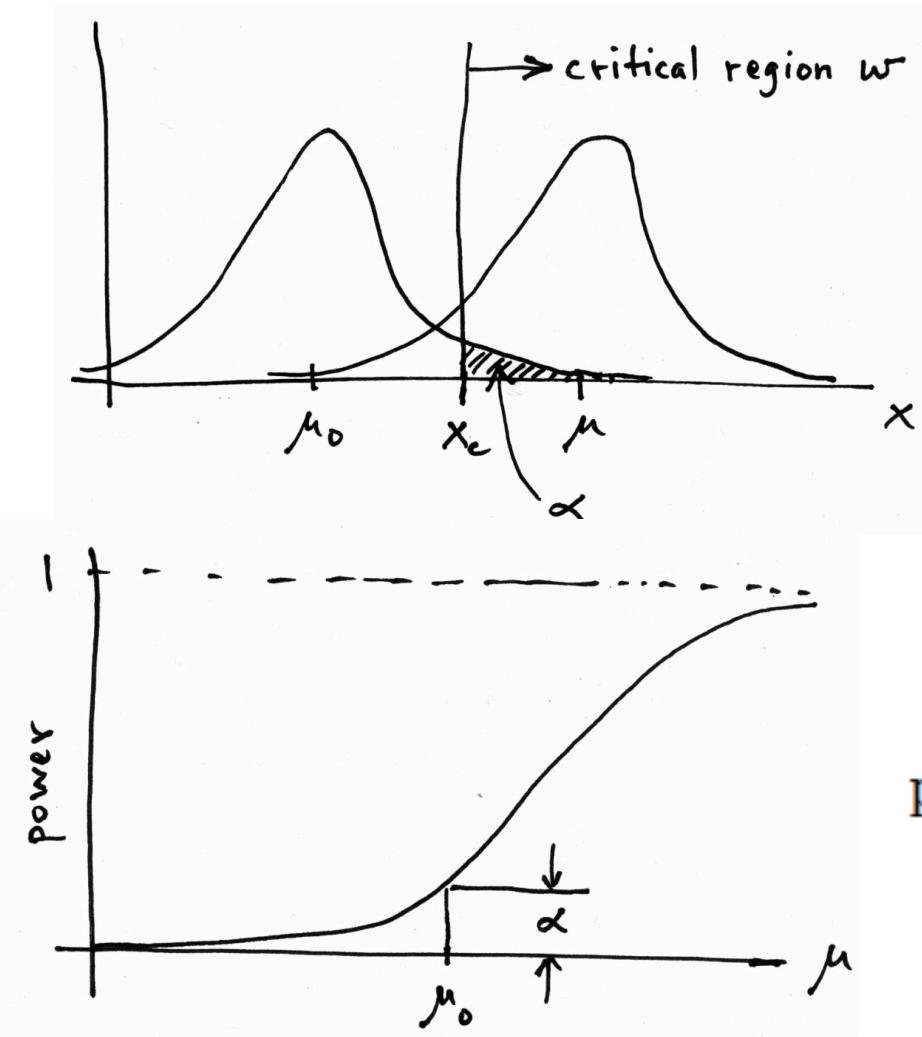
Often such a test has a high power not only with respect to a specific point alternative but for a class of alternatives.  
E.g., using a measurement  $x \sim \text{Gauss}(\mu, \sigma)$  we may test

$H_0 : \mu = \mu_0$  versus the composite alternative  $H_1 : \mu > \mu_0$

We get the highest power with respect to any  $\mu > \mu_0$  by taking the critical region  $x \geq x_c$  where the cut-off  $x_c$  is determined by the significance level such that

$$\alpha = P(x \geq x_c | \mu_0).$$

# Test of $\mu = \mu_0$ vs. $\mu > \mu_0$ with $x \sim \text{Gauss}(\mu, \sigma)$



Standard Gaussian  
cumulative distribution

$$\alpha = 1 - \Phi\left(\frac{x_c - \mu_0}{\sigma}\right)$$

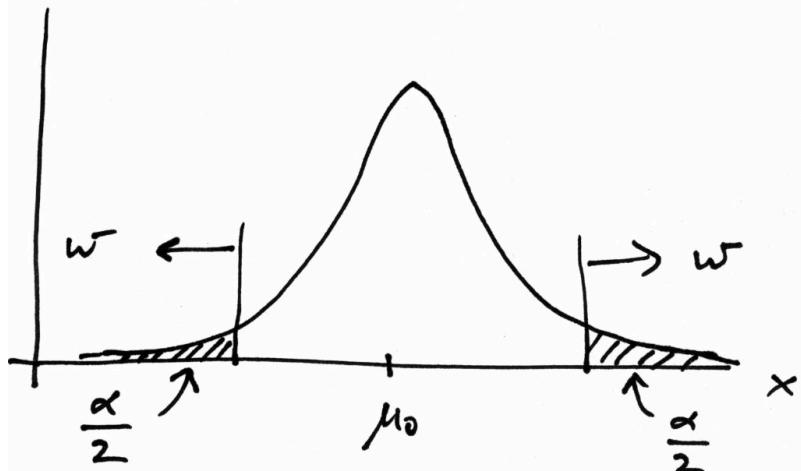
$$x_c = \mu_0 + \sigma\Phi^{-1}(1 - \alpha)$$

Standard Gaussian quantile

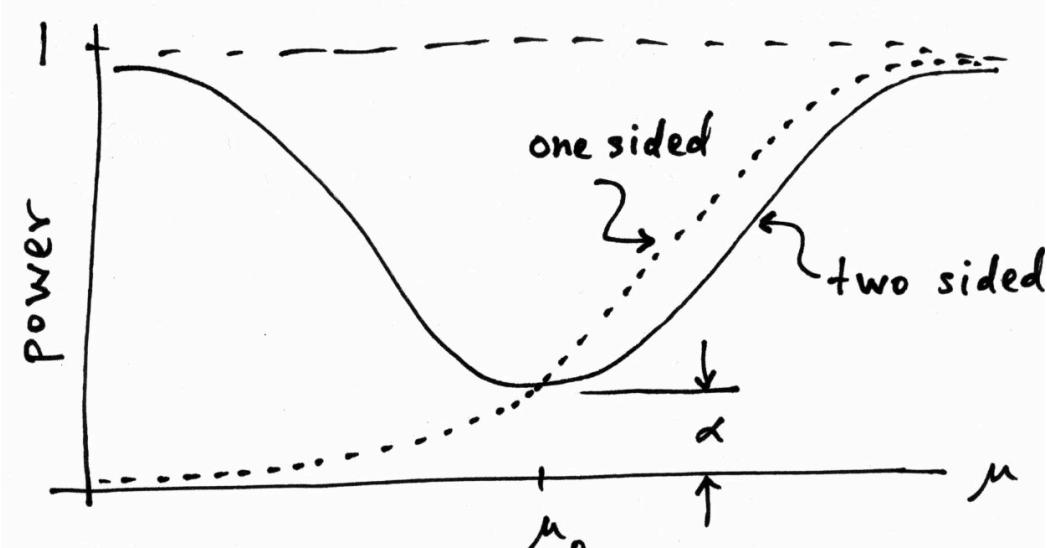
$$\text{power} = 1 - \beta = P(x > x_c | \mu) =$$

$$1 - \Phi\left(\frac{\mu_0 - \mu}{\sigma} + \Phi^{-1}(1 - \alpha)\right)$$

# Choice of critical region based on power (3)



But we might consider  $\mu < \mu_0$  as well as  $\mu > \mu_0$  to be viable alternatives, and choose the critical region to contain both high and low  $x$  (a two-sided test).



New critical region now gives reasonable power for  $\mu < \mu_0$ , but less power for  $\mu > \mu_0$  than the original one-sided test.

# No such thing as a model-independent test

In general we cannot find a single critical region that gives the maximum power for all possible alternatives (no “Uniformly Most Powerful” test).

In HEP we often try to construct a test of

$H_0$  : Standard Model (or “background only”, etc.)

such that we have a well specified “false discovery rate”,

$\alpha$  = Probability to reject  $H_0$  if it is true,

and high power with respect to some interesting alternative,

$H_1$  : SUSY, Z', etc.

But there is no such thing as a “model independent” test. Any statistical test will inevitably have high power with respect to some alternatives and less power with respect to others.

# *p*-values

Suppose hypothesis  $H$  predicts pdf  $f(\vec{x}|H)$  for a set of observations  $\vec{x} = (x_1, \dots, x_n)$ .

We observe a single point in this space:  $\vec{x}_{\text{obs}}$

What can we say about the validity of  $H$  in light of the data?

Express level of compatibility by giving the *p*-value for  $H$ :

$p$  = probability, under assumption of  $H$ , to observe data with equal or lesser compatibility with  $H$  relative to the data we got.

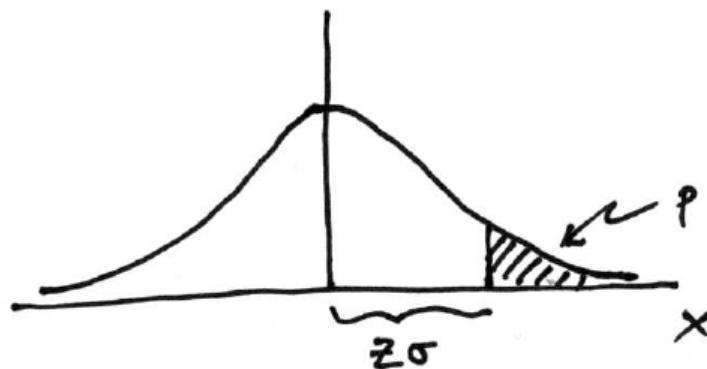


This is not the probability that  $H$  is true!

Requires one to say what part of data space constitutes lesser compatibility with  $H$  than the observed data (implicitly this means that region gives better agreement with some alternative).

# Significance from $p$ -value

Often define significance  $Z$  as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same  $p$ -value.



$$p = \int_Z^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 1 - \Phi(Z) \quad 1 - \text{TMath::Freq}$$

$$Z = \Phi^{-1}(1 - p) \quad \text{TMath::NormQuantile}$$

E.g.  $Z = 5$  (a “5 sigma effect”) corresponds to  $p = 2.9 \times 10^{-7}$ .

# Using a $p$ -value to define test of $H_0$

One can show the distribution of the  $p$ -value of  $H$ , under assumption of  $H$ , is uniform in  $[0,1]$ .

So the probability to find the  $p$ -value of  $H_0$ ,  $p_0$ , less than  $\alpha$  is

$$P(p_0 \leq \alpha | H_0) = \alpha$$

We can define the critical region of a test of  $H_0$  with size  $\alpha$  as the set of data space where  $p_0 \leq \alpha$ .

Formally the  $p$ -value relates only to  $H_0$ , but the resulting test will have a given power with respect to a given alternative  $H_1$ .

# The Poisson counting experiment

Suppose we do a counting experiment and observe  $n$  events.

Events could be from *signal* process or from *background* – we only count the total number.

Poisson model:

$$P(n|s, b) = \frac{(s + b)^n}{n!} e^{-(s+b)}$$

$s$  = mean (i.e., expected) # of signal events

$b$  = mean # of background events

Goal is to make inference about  $s$ , e.g.,

test  $s = 0$  (rejecting  $H_0 \approx$  “discovery of signal process”)

test all non-zero  $s$  (values not rejected = confidence interval)

In both cases need to ask what is relevant alternative hypothesis.

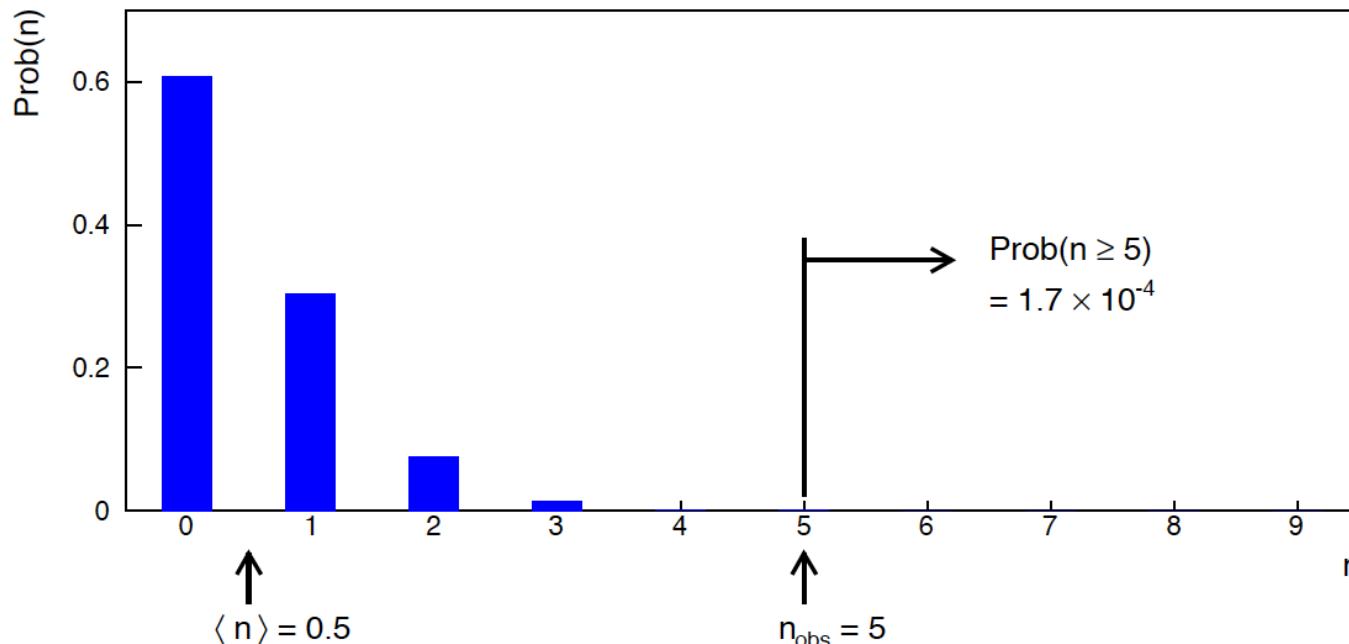
# Poisson counting experiment: discovery $p$ -value

Suppose  $b = 0.5$  (known), and we observe  $n_{\text{obs}} = 5$ .

Should we claim evidence for a new discovery?

Give  $p$ -value for hypothesis  $s = 0$ :

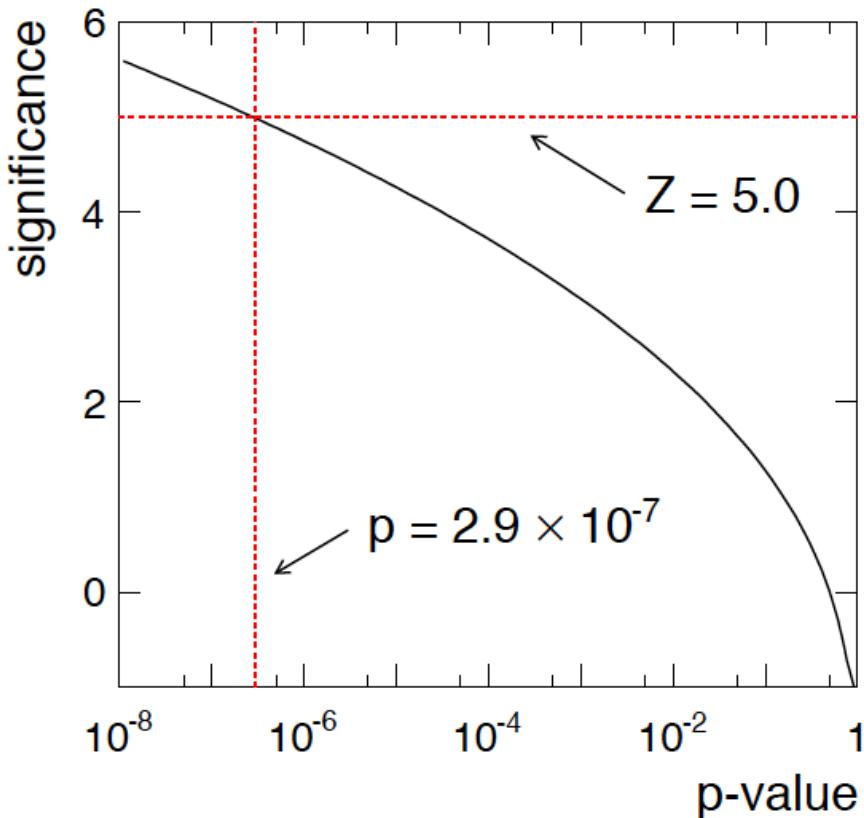
$$\begin{aligned} p\text{-value} &= P(n \geq 5; b = 0.5, s = 0) \\ &= 1.7 \times 10^{-4} \neq P(s = 0)! \end{aligned}$$



# Poisson counting experiment: discovery significance

Equivalent significance for  $p = 1.7 \times 10^{-4}$ :  $Z = \Phi^{-1}(1 - p) = 3.6$

Often claim discovery if  $Z > 5$  ( $p < 2.9 \times 10^{-7}$ , i.e., a “5-sigma effect”)



In fact this tradition should be revisited:  $p$ -value intended to quantify probability of a signal-like fluctuation assuming background only; not intended to cover, e.g., hidden systematics, plausibility signal model, compatibility of data with signal, “look-elsewhere effect” (~multiple testing), etc.

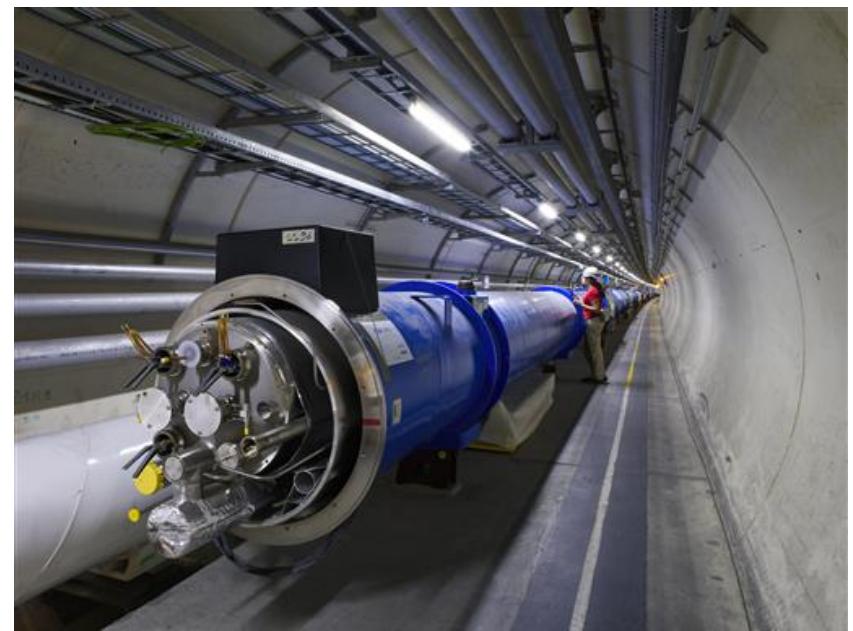
# Example setting for statistical tests: the Large Hadron Collider



Counter-rotating proton beams  
in 27 km circumference ring

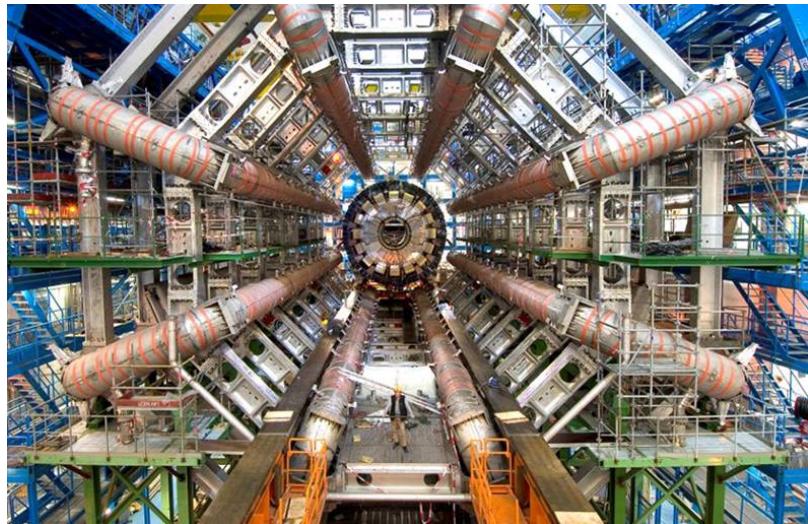
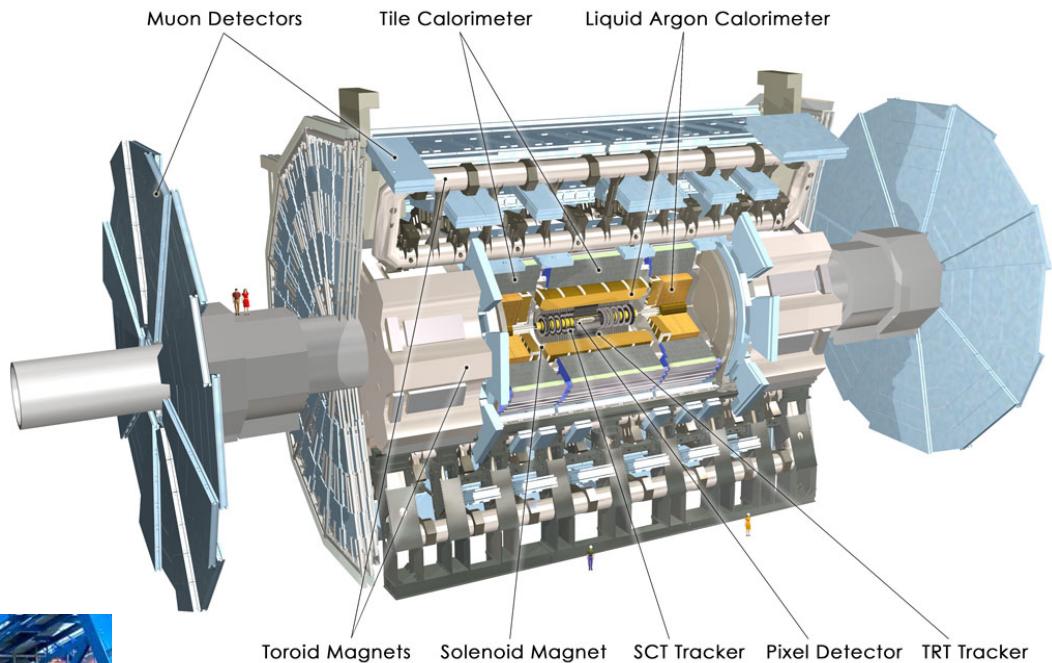
pp centre-of-mass energy 14 TeV

Detectors at 4 pp collision points:  
ATLAS  
CMS ← general purpose  
LHCb (b physics)  
ALICE (heavy ion physics)



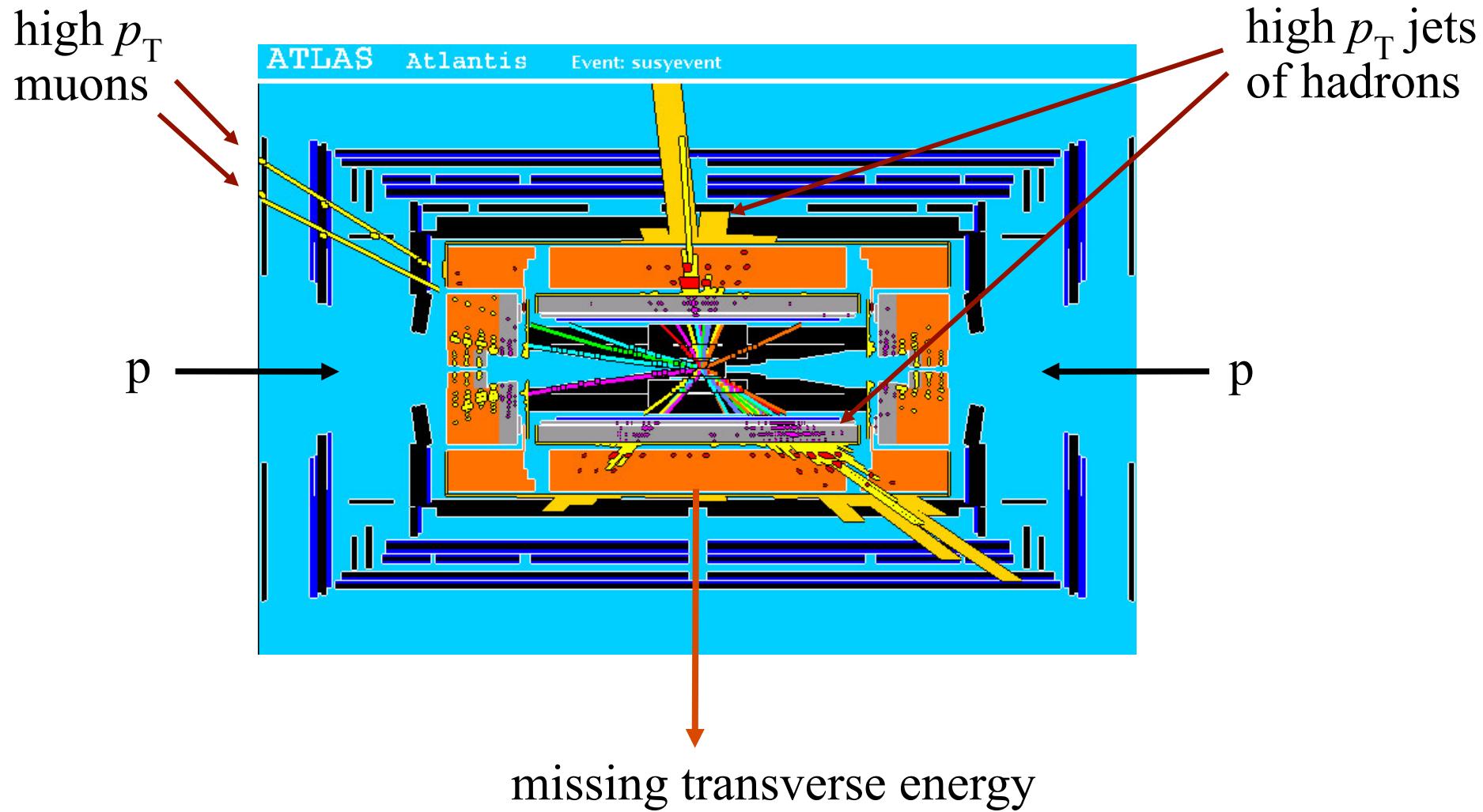
# The ATLAS detector

2100 physicists  
37 countries  
167 universities/labs

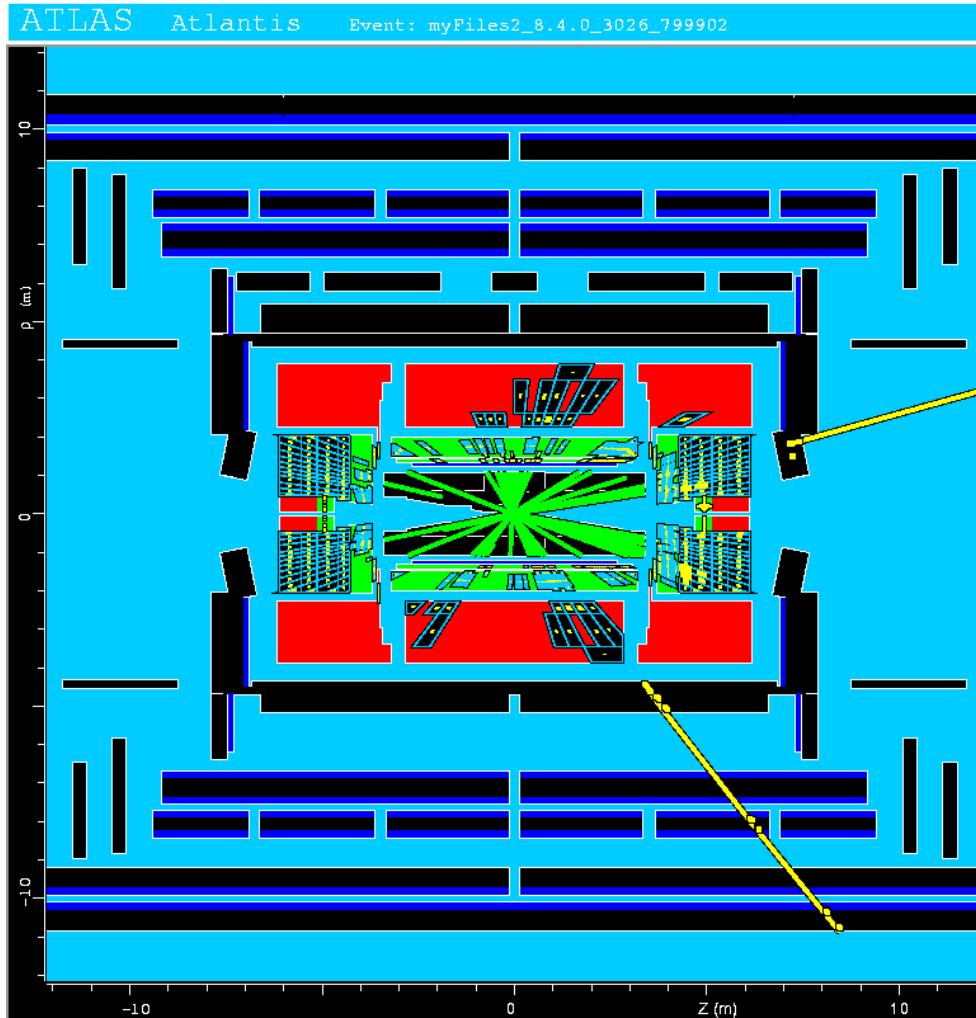


25 m diameter  
46 m length  
7000 tonnes  
 $\sim 10^8$  electronic channels

# A simulated SUSY event



# Background events



This event from Standard Model  $t\bar{t}$  production also has high  $p_T$  jets and muons, and some missing transverse energy.

→ can easily mimic a SUSY event.

# Statistical tests (in a particle physics context)

Suppose the result of a measurement for an individual event is a collection of numbers  $\vec{x} = (x_1, \dots, x_n)$

$x_1$  = number of muons,

$x_2$  = mean  $p_T$  of jets,

$x_3$  = missing energy, ...

$\vec{x}$  follows some  $n$ -dimensional joint pdf, which depends on the type of event produced, i.e., was it

$$pp \rightarrow t\bar{t}, \quad pp \rightarrow \tilde{g}\tilde{g}, \dots$$

For each reaction we consider we will have a hypothesis for the pdf of  $\vec{x}$ , e.g.,  $f(\vec{x}|H_0)$ ,  $f(\vec{x}|H_1)$ , etc.

E.g. call  $H_0$  the **background** hypothesis (the event type we want to reject);  $H_1$  is **signal** hypothesis (the type we want).

# Selecting events

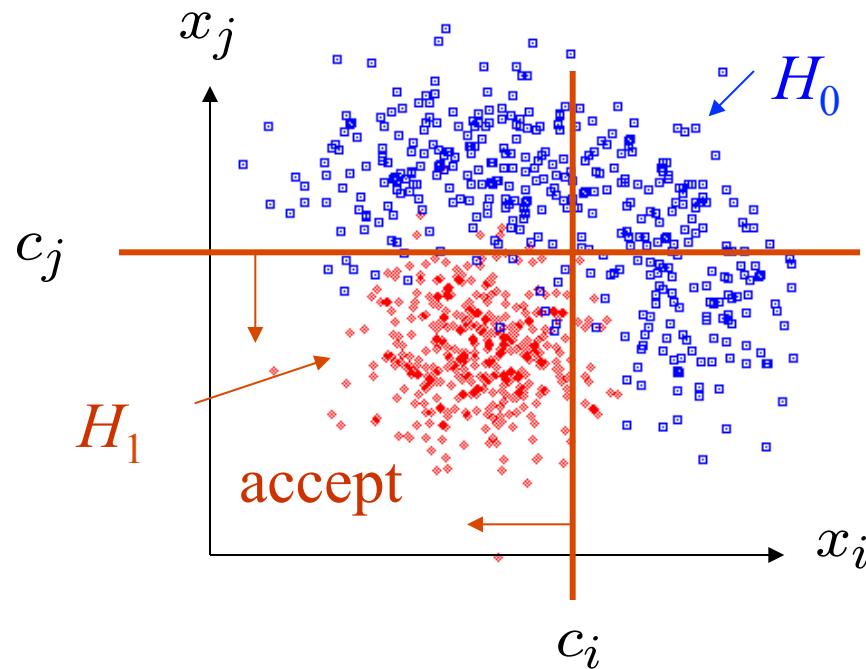
Suppose we have a data sample with two kinds of events, corresponding to hypotheses  $H_0$  and  $H_1$  and we want to select those of type  $H_1$ .

Each event is a point in  $\vec{x}$  space. What ‘decision boundary’ should we use to accept/reject events as belonging to event types  $H_0$  or  $H_1$ ?

Perhaps select events with ‘cuts’:

$$x_i < c_i$$

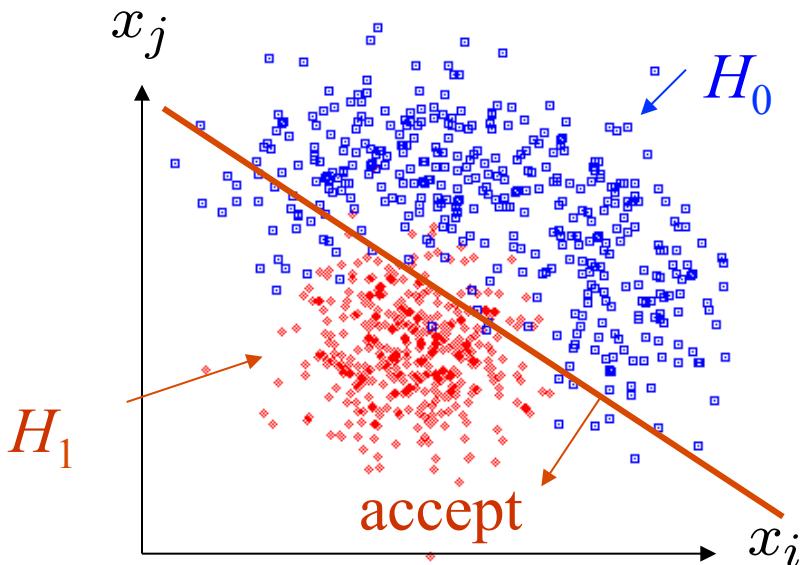
$$x_j < c_j$$



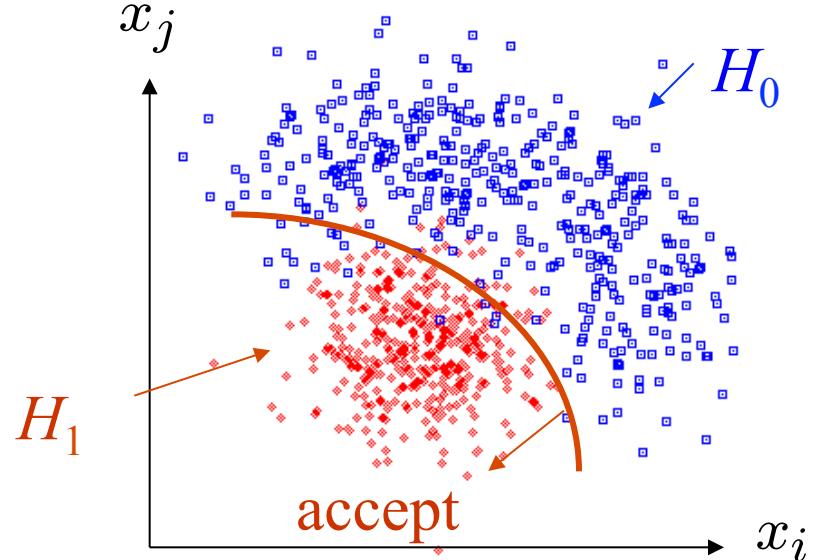
# Other ways to select events

Or maybe use some other sort of decision boundary:

linear



or nonlinear



How can we do this in an ‘optimal’ way?

# Test statistics

The boundary of the critical region for an  $n$ -dimensional data space  $\mathbf{x} = (x_1, \dots, x_n)$  can be defined by an equation of the form

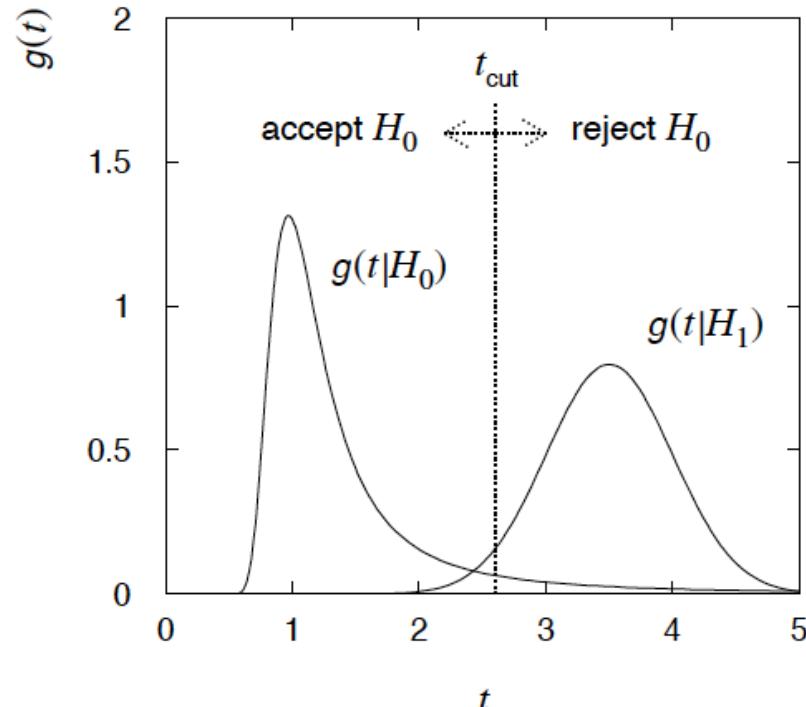
$$t(x_1, \dots, x_n) = t_{\text{cut}}$$

where  $t(x_1, \dots, x_n)$  is a scalar test statistic.

We can work out the pdfs  $g(t|H_0)$ ,  $g(t|H_1)$ , ...

Decision boundary is now a single ‘cut’ on  $t$ , defining the critical region.

So for an  $n$ -dimensional problem we have a corresponding 1-d problem.



# Test statistic based on likelihood ratio

How can we choose a test's critical region in an ‘optimal way’?

Neyman-Pearson lemma states:

To get the highest power for a given significance level in a test of  $H_0$ , (background) versus  $H_1$ , (signal) the critical region should have

$$\frac{f(\mathbf{x}|H_1)}{f(\mathbf{x}|H_0)} > c$$

inside the region, and  $\leq c$  outside, where  $c$  is a constant chosen to give a test of the desired size.

Equivalently, optimal scalar test statistic is

$$t(\mathbf{x}) = \frac{f(\mathbf{x}|H_1)}{f(\mathbf{x}|H_0)}$$

N.B. any monotonic function of this leads to the same test.

# Classification viewed as a statistical test

Probability to reject  $H_0$  if true (type I error):  $\alpha = \int_W f(\mathbf{x}|H_0) d\mathbf{x}$

$\alpha$  = size of test, significance level, false discovery rate

Probability to accept  $H_0$  if  $H_1$  true (type II error)  $\beta = \int_{\overline{W}} f(\mathbf{x}|H_1) d\mathbf{x}$

$\beta$  = power of test with respect to  $H_1$

Equivalently if e.g.  $H_0$  = background,  $H_1$  = signal, use efficiencies:

$$\varepsilon_b = \int_W f(\mathbf{x}|H_0) = \alpha$$

$$\varepsilon_s = \int_W f(\mathbf{x}|H_1) = 1 - \beta = \text{power}$$

# Purity / misclassification rate

Consider the probability that an event of signal (s) type classified correctly (i.e., the event selection purity),

Use Bayes' theorem:

Here  $W$  is signal region

$$P(s|x \in W) = \frac{P(x \in W|s)P(s)}{P(x \in W|s)P(s) + P(x \in W|b)P(b)}$$

posterior probability = signal purity  
=  $1 -$  signal misclassification rate

Note purity depends on the prior probability for an event to be signal or background as well as on s/b efficiencies.

# Neyman-Pearson doesn't usually help

We usually don't have explicit formulae for the pdfs  $f(\mathbf{x}|s)$ ,  $f(\mathbf{x}|b)$ , so for a given  $\mathbf{x}$  we can't evaluate the likelihood ratio

$$t(\mathbf{x}) = \frac{f(\mathbf{x}|s)}{f(\mathbf{x}|b)}$$

Instead we may have Monte Carlo models for signal and background processes, so we can produce simulated data:

generate  $\mathbf{x} \sim f(\mathbf{x}|s) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_N$

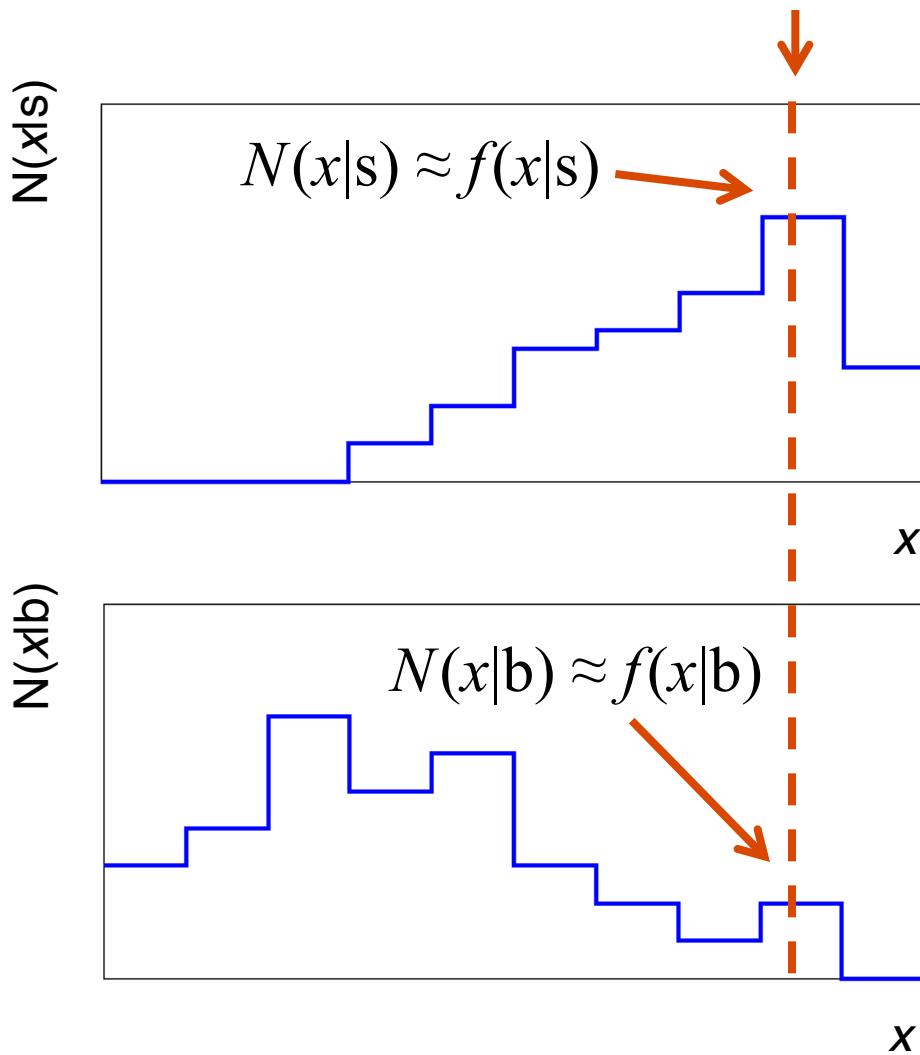
generate  $\mathbf{x} \sim f(\mathbf{x}|b) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_N$

This gives samples of “training data” with events of known type.

Can be expensive (1 fully simulated LHC event  $\sim 1$  CPU minute).

# Approximate LR from histograms

Want  $t(x) = f(x|s)/f(x|b)$  for  $x$  here



One possibility is to generate MC data and construct histograms for both signal and background.

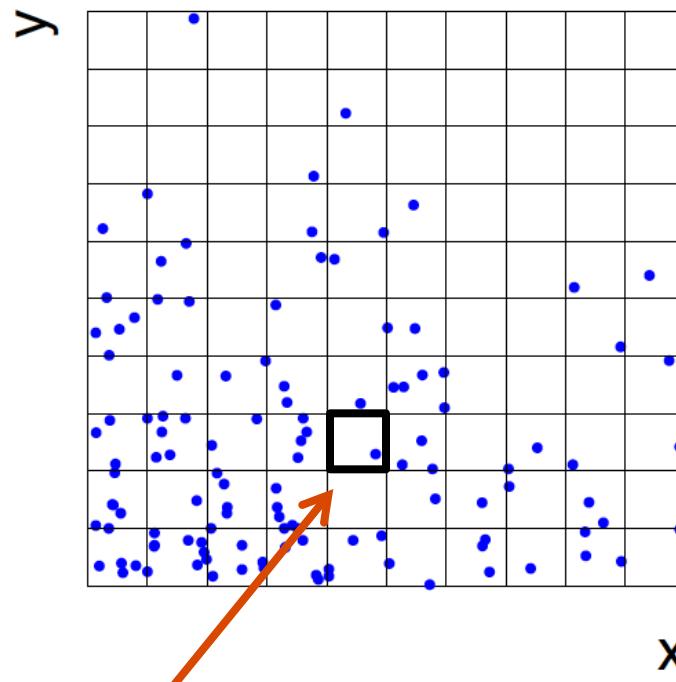
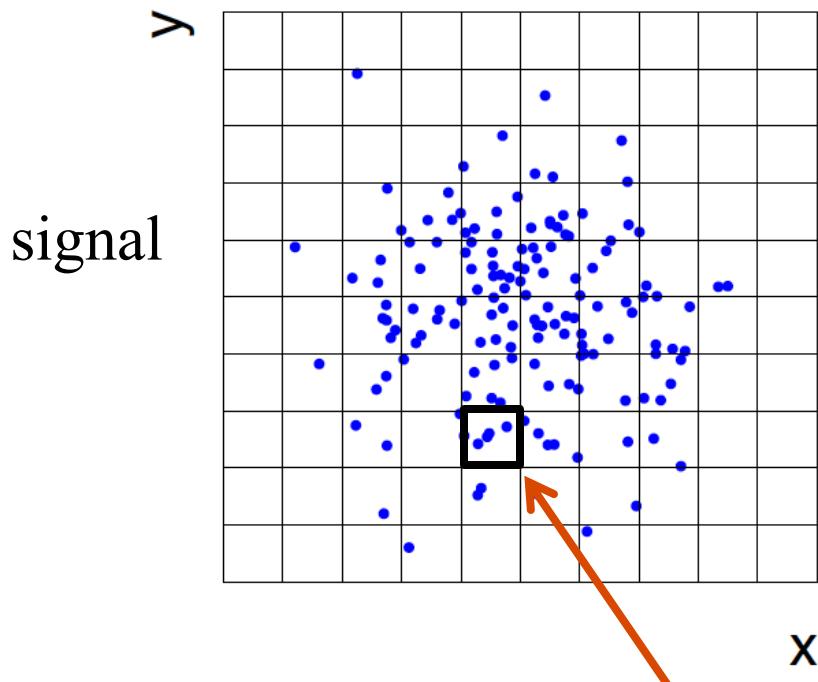
Use (normalized) histogram values to approximate LR:

$$t(x) \approx \frac{N(x|s)}{N(x|b)}$$

Can work well for single variable.

# Approximate LR from 2D-histograms

Suppose problem has 2 variables. Try using 2-D histograms:



Approximate pdfs using  $N(x,y|s)$ ,  $N(x,y|b)$  in corresponding cells.

But if we want  $M$  bins for each variable, then in  $n$ -dimensions we have  $M^n$  cells; can't generate enough training data to populate.

→ Histogram method usually not usable for  $n > 1$  dimension.

# Strategies for multivariate analysis

Neyman-Pearson lemma gives optimal answer, but cannot be used directly, because we usually don't have  $f(\mathbf{x}|s), f(\mathbf{x}|b)$ .

Histogram method with  $M$  bins for  $n$  variables requires that we estimate  $M^n$  parameters (the values of the pdfs in each cell), so this is rarely practical.

A compromise solution is to assume a certain functional form for the test statistic  $t(\mathbf{x})$  with fewer parameters; determine them (using MC) to give best separation between signal and background.

Alternatively, try to estimate the probability densities  $f(\mathbf{x}|s)$  and  $f(\mathbf{x}|b)$  (with something better than histograms) and use the estimated pdfs to construct an approximate likelihood ratio.

# Multivariate methods

Many new (and some old) methods:

Fisher discriminant

Neural networks

Kernel density methods

Support Vector Machines

Decision trees

Boosting

Bagging

New software for HEP, e.g.,

**TMVA** , Höcker, Stelzer, Tegenfeldt, Voss, Voss, physics/0703039

**StatPatternRecognition**, I. Narsky, physics/0507143

# Resources on multivariate methods

C.M. Bishop, Pattern Recognition and Machine Learning,  
Springer, 2006

T. Hastie, R. Tibshirani, J. Friedman, The Elements of  
Statistical Learning, 2<sup>nd</sup> ed., Springer, 2009

R. Duda, P. Hart, D. Stork, Pattern Classification, 2<sup>nd</sup> ed.,  
Wiley, 2001

A. Webb, Statistical Pattern Recognition, 2<sup>nd</sup> ed., Wiley, 2002.

Ilya Narsky and Frank C. Porter, *Statistical Analysis  
Techniques in Particle Physics*, Wiley, 2014.

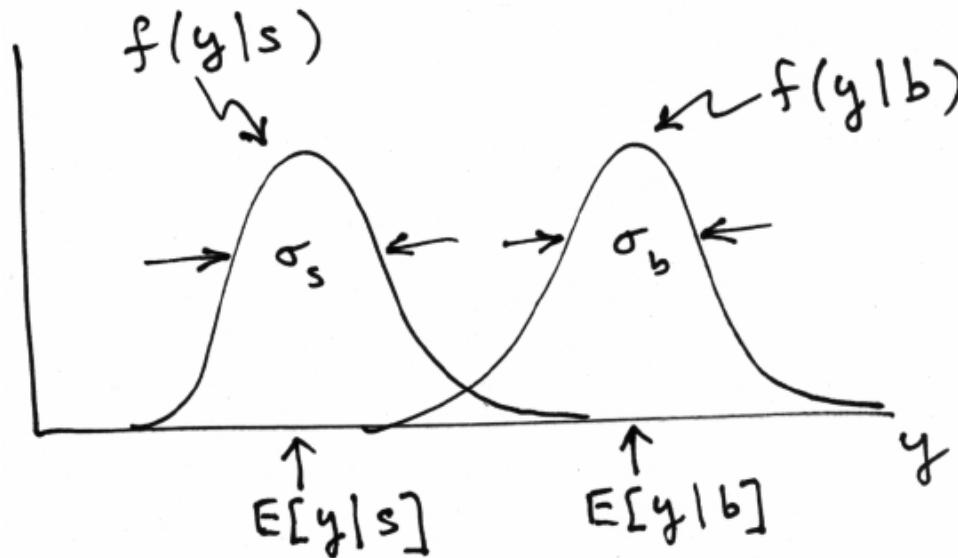
朱永生（编著），实验数据多元统计分析，科学出版社，  
北京，2009。

# Linear test statistic

Suppose there are  $n$  input variables:  $\mathbf{x} = (x_1, \dots, x_n)$ .

Consider a linear function:  $y(\mathbf{x}) = \sum_{i=1}^n w_i x_i$

For a given choice of the coefficients  $\mathbf{w} = (w_1, \dots, w_n)$  we will get pdfs  $f(y|s)$  and  $f(y|b)$ :



# Linear test statistic

Fisher: to get large difference between means and small widths for  $f(y|s)$  and  $f(y|b)$ , maximize the difference squared of the expectation values divided by the sum of the variances:

$$J(\mathbf{w}) = \frac{(E[y|s] - E[y|b])^2}{V[y|s] + V[y|b]}$$

Setting  $\partial J / \partial w_i = 0$  gives:

$$\mathbf{w} \propto W^{-1}(\boldsymbol{\mu}_b - \boldsymbol{\mu}_s)$$

$$W_{ij} = \text{cov}[x_i, x_j|s] + \text{cov}[x_i, x_j|b]$$

$$\mu_{i,s} = E[x_i|s], \quad \mu_{i,b} = E[x_i|b]$$

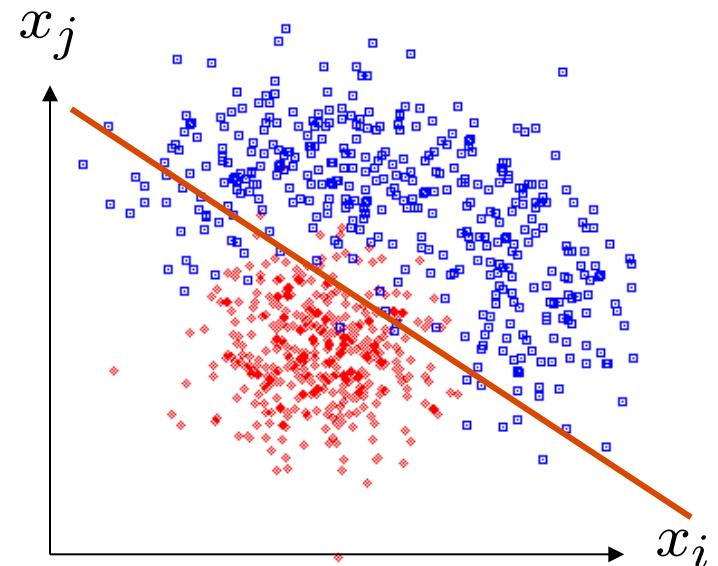
# The Fisher discriminant

The resulting coefficients  $w_i$  define a Fisher discriminant.

Coefficients defined up to multiplicative constant; can also add arbitrary offset, i.e., usually define test statistic as

$$y(\mathbf{x}) = w_0 + \sum_{i=1}^n w_i x_i$$

Boundaries of the test's critical region are surfaces of constant  $y(\mathbf{x})$ , here linear (hyperplanes):

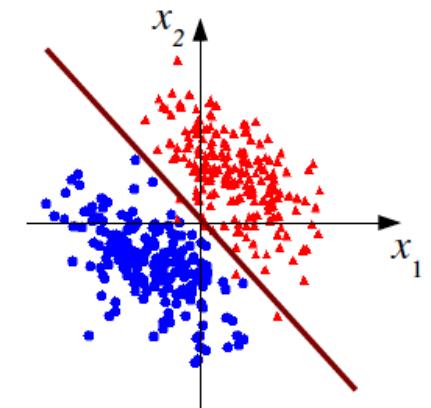


# Fisher discriminant for Gaussian data

Suppose the pdfs of the input variables,  $f(\mathbf{x}|s)$  and  $f(\mathbf{x}|b)$ , are both multivariate Gaussians with same covariance but different means:

$$\begin{aligned}f(\mathbf{x}|s) &= \text{Gauss}(\boldsymbol{\mu}_s, V) \\f(\mathbf{x}|b) &= \text{Gauss}(\boldsymbol{\mu}_b, V)\end{aligned}$$

Same covariance  
 $V_{ij} = \text{cov}[x_i, x_j]$



In this case it can be shown that the Fisher discriminant is

$$y(\mathbf{x}) \sim \ln \frac{f(\mathbf{x}|s)}{f(\mathbf{x}|b)}$$

i.e., it is a monotonic function of the likelihood ratio and thus leads to the same critical region. So in this case the Fisher discriminant provides an optimal statistical test.

# Transformation of inputs

If the data are not Gaussian with equal covariance, a linear decision boundary is not optimal. But we can try to subject the data to a transformation

$$\varphi_1(\vec{x}), \dots, \varphi_m(\vec{x})$$

and then treat the  $\varphi_i$  as the new input variables. This is often called “feature space” and the  $\varphi_i$  are “basis functions”. The basis functions can be fixed or can contain adjustable parameters which we optimize with training data (cf. neural networks).

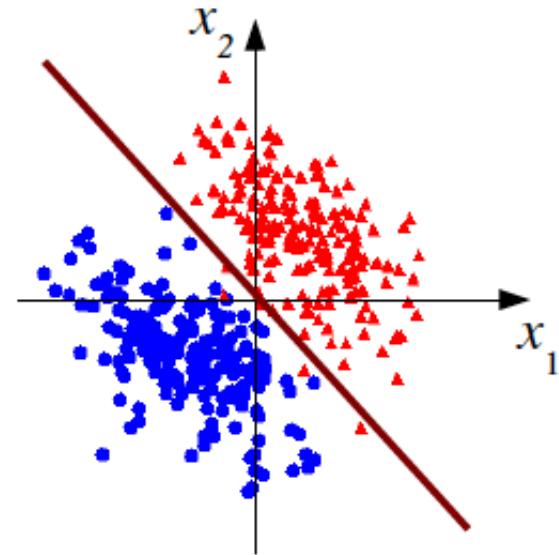
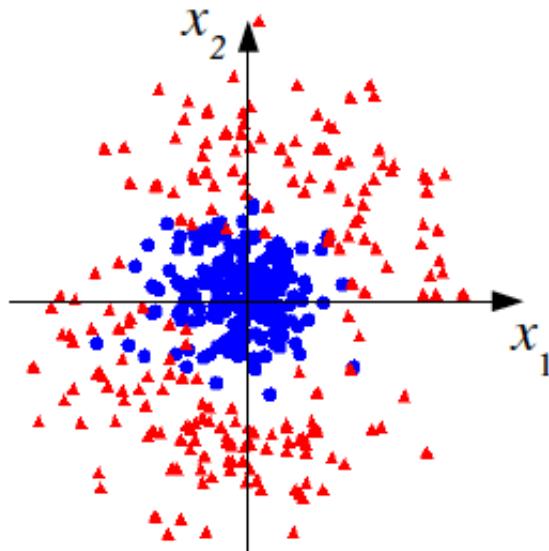
In other cases we will see that the basis functions only enter as dot products

$$\vec{\varphi}(\vec{x}_i) \cdot \vec{\varphi}(\vec{x}_j) = K(\vec{x}_i, \vec{x}_j)$$

and thus we will only need the “kernel function”  $K(x_i, x_j)$

# Linear decision boundaries

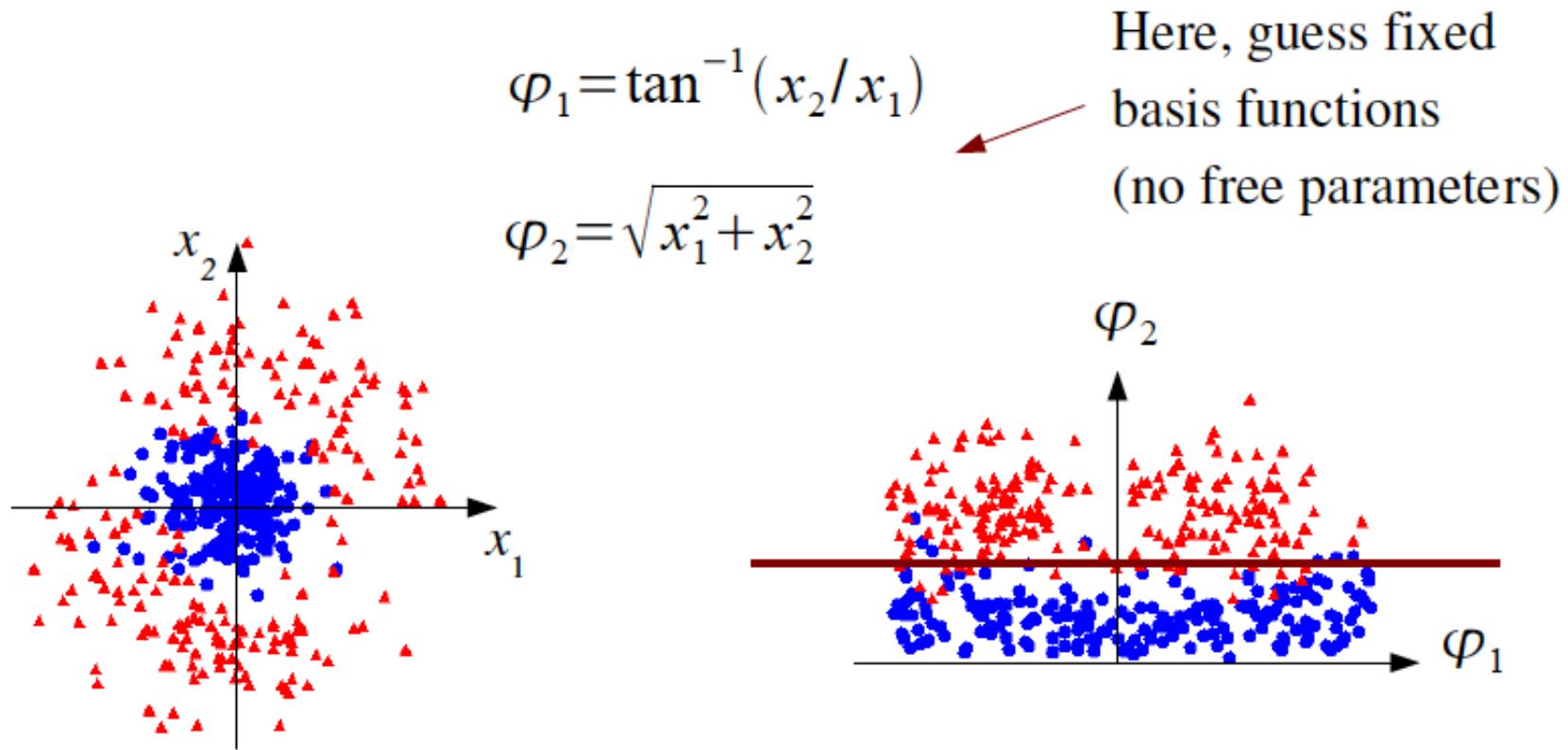
A linear decision boundary is only optimal when both classes follow multivariate Gaussians with equal covariances and different means.



For some other cases a linear boundary is almost useless.

# Nonlinear transformation of inputs

We can try to find a transformation,  $x_1, \dots, x_n \rightarrow \varphi_1(\vec{x}), \dots, \varphi_m(\vec{x})$  so that the transformed “feature space” variables can be separated better by a linear boundary:



# Neural networks

Neural networks originate from attempts to model neural processes (McCulloch and Pitts, 1943; Rosenblatt, 1962).

Widely used in many fields, and for many years the only “advanced” multivariate method popular in HEP.

We can view a neural network as a specific way of parametrizing the basis functions used to define the feature space transformation.

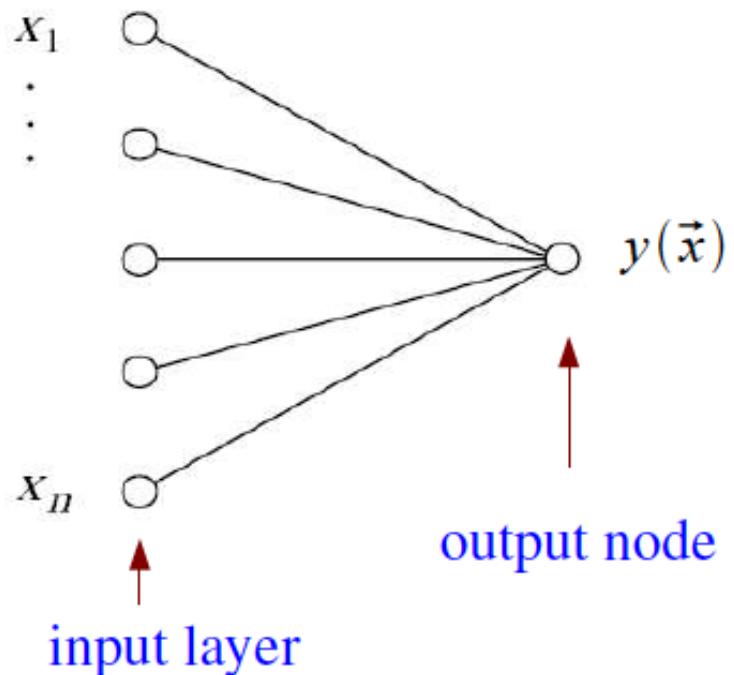
The training data are then used to adjust the parameters so that the resulting discriminant function has the best performance.

# The single layer perceptron

Define the discriminant using  $y(\vec{x}) = h\left(w_0 + \sum_{i=1}^n w_i x_i\right)$

where  $h$  is a nonlinear, monotonic activation function; we can use e.g. the logistic sigmoid  $h(x) = (1 + e^{-x})^{-1}$ .

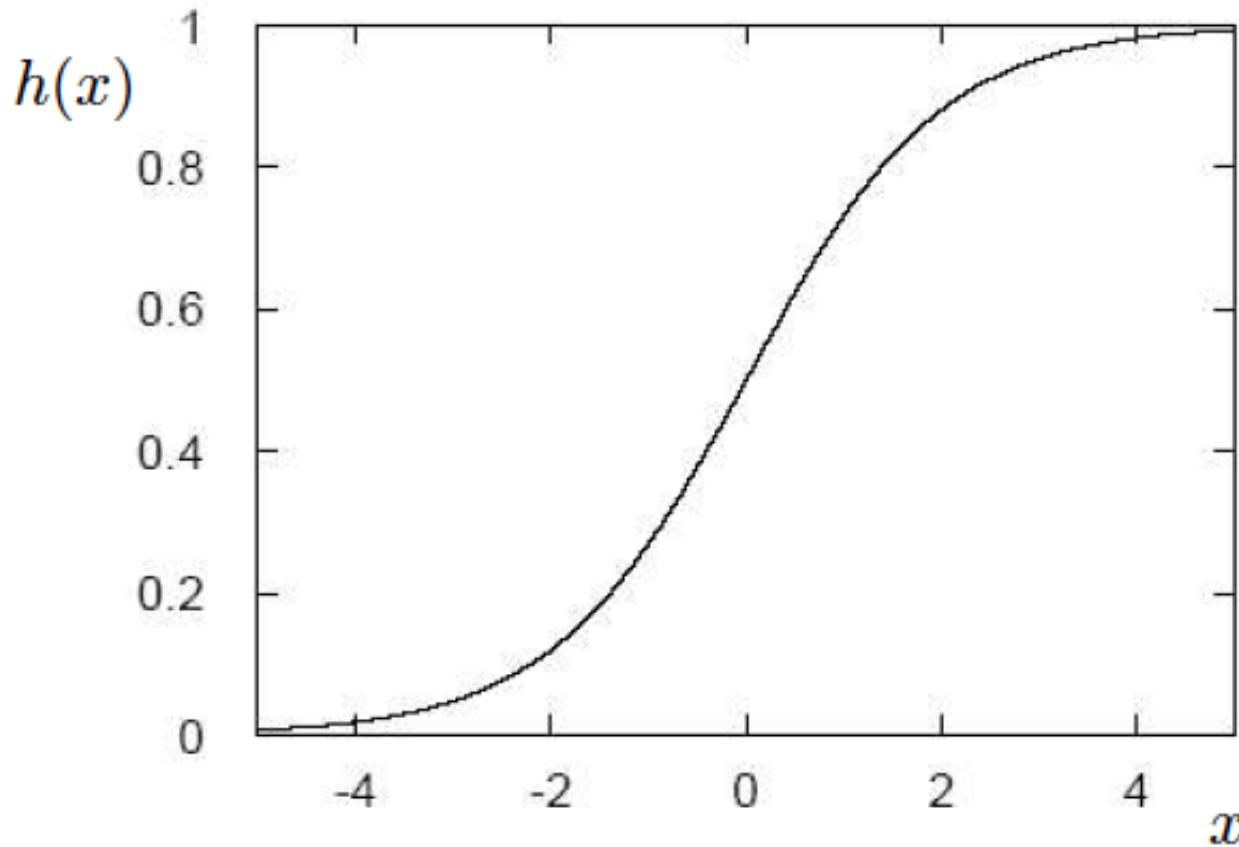
If the activation function is monotonic, the resulting  $y(\vec{x})$  is equivalent to the original linear discriminant. This is an example of a “generalized linear model” called the **single layer perceptron**.



# The activation function

For activation function  $h(\cdot)$  often use logistic sigmoid:

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



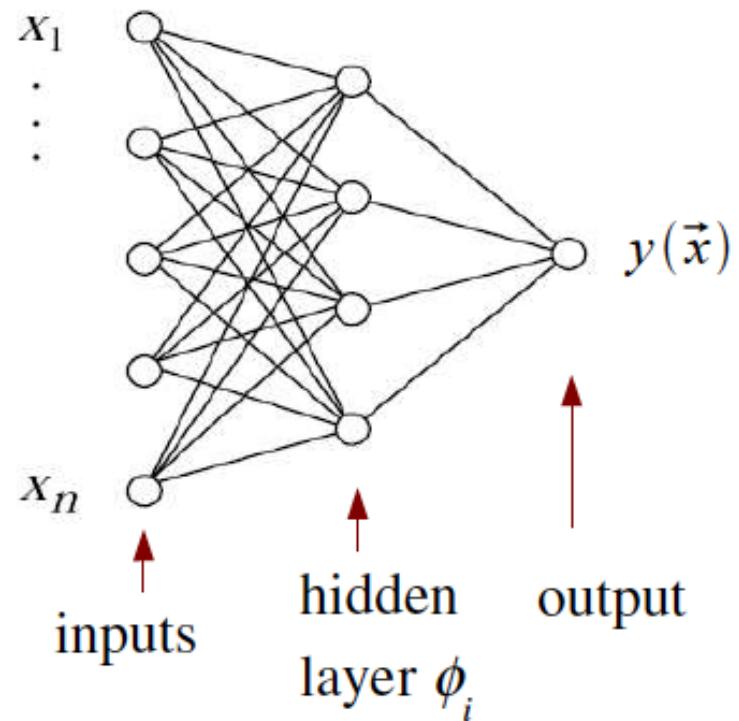
# The multilayer perceptron

Now use this idea to define not only the output  $y(\vec{x})$ , but also the set of transformed inputs  $\varphi_1(\vec{x}), \dots, \varphi_m(\vec{x})$  that form a “hidden layer”:

Superscript for weights indicates  
layer number

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

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



This is the **multilayer perceptron**, our basic neural network model;  
straightforward to generalize to multiple hidden layers.

# Network architecture: one hidden layer

**Theorem:** An MLP with a single hidden layer having a sufficiently large number of nodes can approximate arbitrarily well the optimal decision boundary.

Holds for any continuous non-polynomial activation function

Leshno, Lin, Pinkus and Schocken (1993), *Neural Networks* 6, 861—867

In practice often choose a single hidden layer and try increasing the the number of nodes until no further improvement in performance is found.

# More than one hidden layer

“Relatively little is known concerning the advantages and disadvantages of using a single hidden layer with many units (neurons) over many hidden layers with fewer units. The mathematics and approximation theory of the MLP model with more than one hidden layer is not 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 in neural networks*,  
Acta Numerica (1999), pp. 143—195.

# Network training

The type of each training event is known, i.e., for event  $a$  we have:

$\vec{x}_a = (x_1, \dots, x_n)$  the input variables, and

$t_a = 0, 1$  a numerical label for event type (“target value”)

Let  $\mathbf{w}$  denote the set of all of the weights of the network. We can determine their optimal values by minimizing a sum-of-squares “error function”

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



Contribution to error function  
from each event

# Numerical minimization of $E(\mathbf{w})$

Consider gradient descent method: from an initial guess in weight space  $\mathbf{w}^{(1)}$  take a small step in the direction of maximum decrease.  
I.e. for the step  $\tau$  to  $\tau+1$ ,

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})$$

↑  
learning rate ( $\eta > 0$ )

If we do this with the full error function  $E(\mathbf{w})$ , gradient descent does surprisingly poorly; better to use “conjugate gradients”.

But gradient descent turns out to be useful with an online (sequential) method, i.e., where we update  $\mathbf{w}$  for each training event  $a$ , (cycle through all training events):

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

# Error backpropagation

Error backpropagation (“backprop”) is an algorithm for finding the derivatives required for gradient descent minimization.

The network output can be written  $y(\vec{x}) = h(u(\vec{x}))$  where

$$u(\vec{x}) = \sum_{j=0} w_{1j}^{(2)} \varphi_j(\vec{x}), \quad \varphi_j(\vec{x}) = h\left(\sum_{k=0} w_{jk}^{(1)} x_k\right)$$

where we defined  $\phi_0 = x_0 = 1$  and wrote the sums over the nodes in the preceding layers starting from 0 to include the offsets.

So e.g. for event  $a$  we have  $\frac{\partial E_a}{\partial w_{1j}^{(2)}} = (y_a - t_a) h'(u(\vec{x})) \varphi_j(\vec{x})$



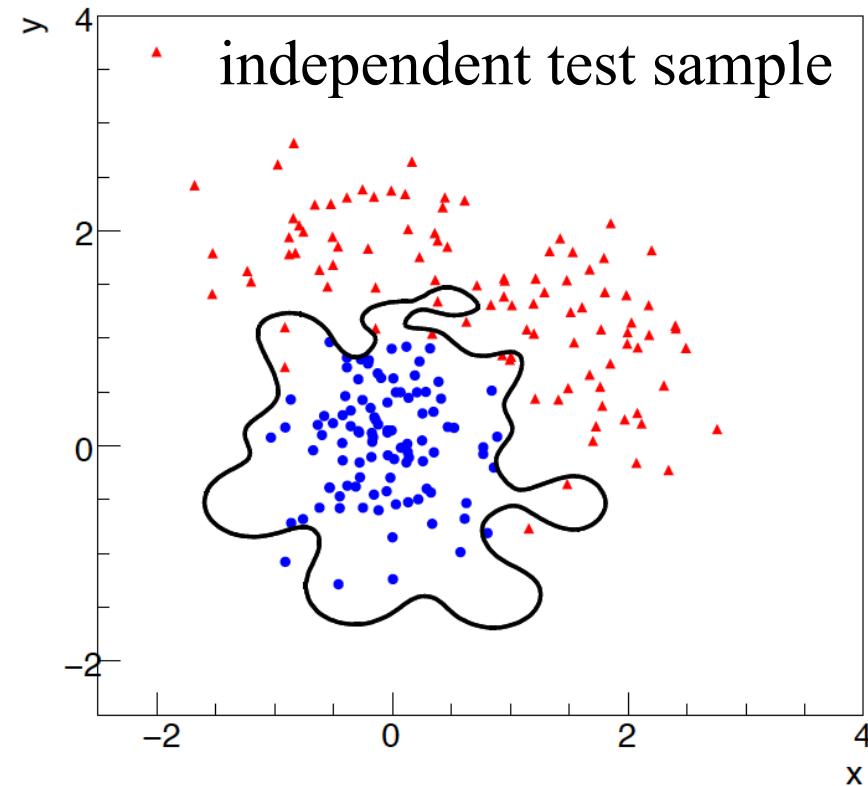
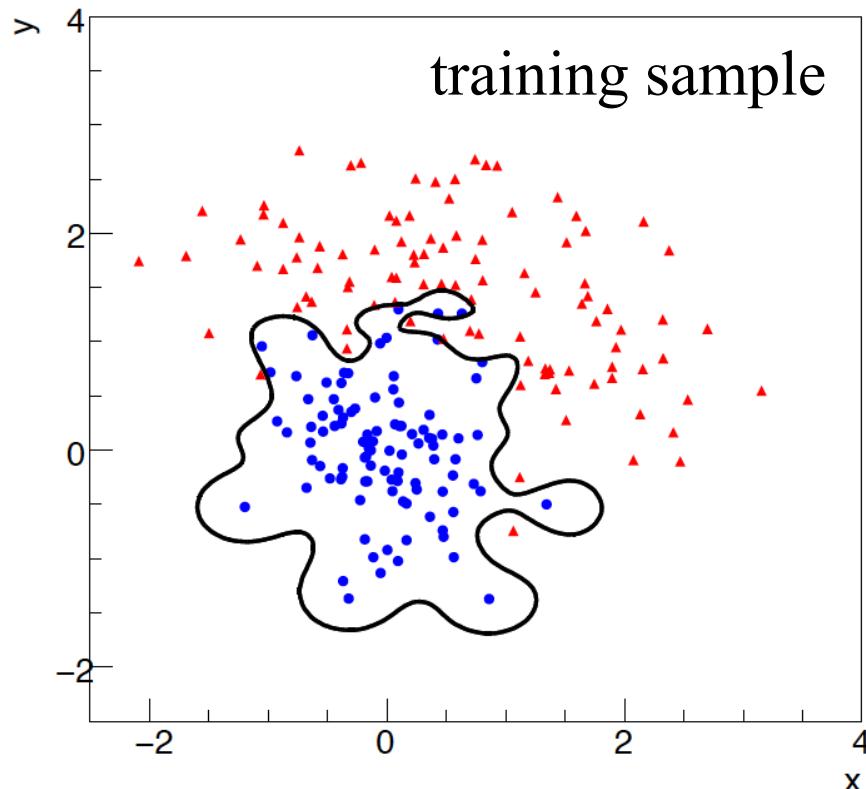
derivative of  
activation function

Chain rule gives all the needed derivatives.

# Overtraining

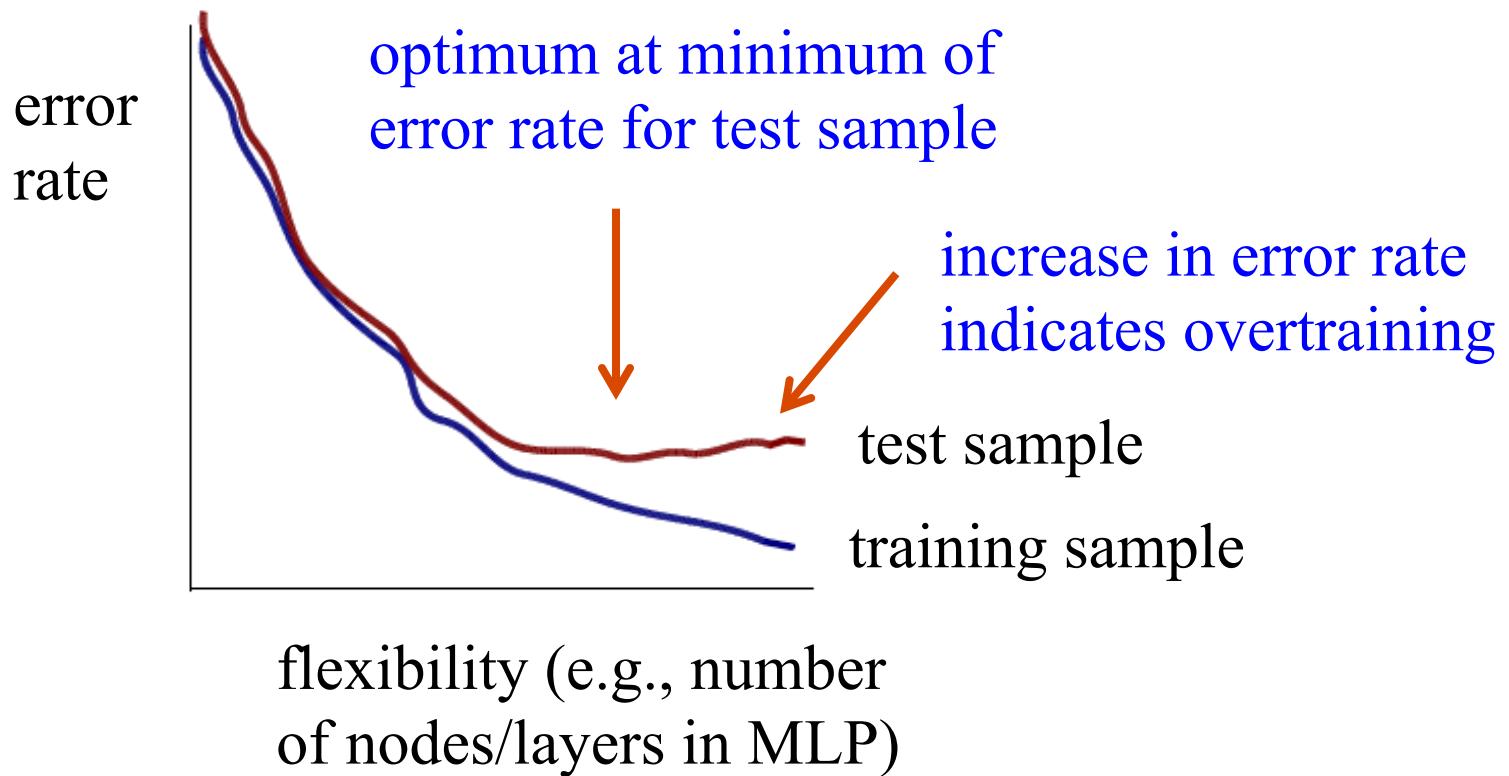
Including more parameters in a classifier makes its decision boundary increasingly flexible, e.g., more nodes/layers for a neural network.

A “flexible” classifier may conform too closely to the training points; the same boundary will not perform well on an independent test data sample ( $\rightarrow$  “overtraining”).



# Monitoring overtraining

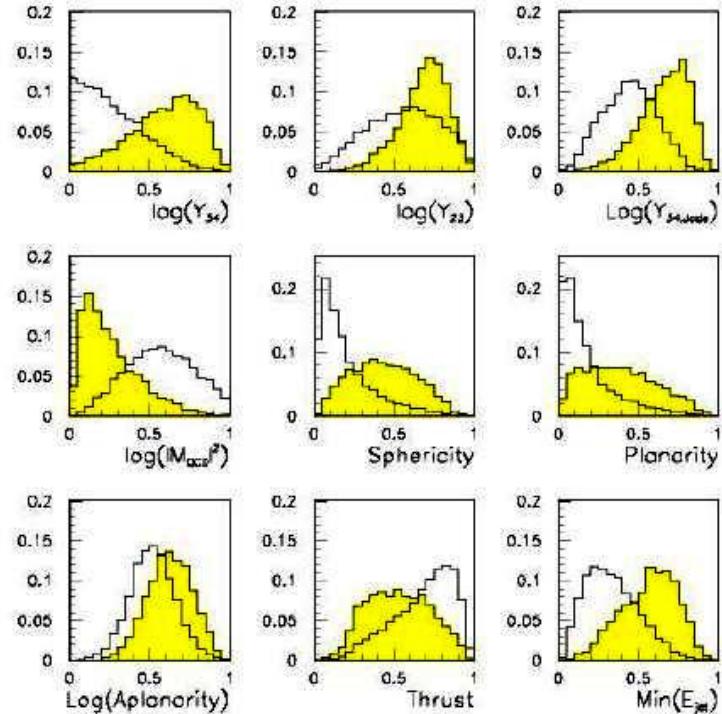
If we monitor the fraction of misclassified events (or similar, e.g., error function  $E(\mathbf{w})$ ) for test and training samples, it will usually decrease for both as the boundary is made more flexible:



# Neural network example from LEP II

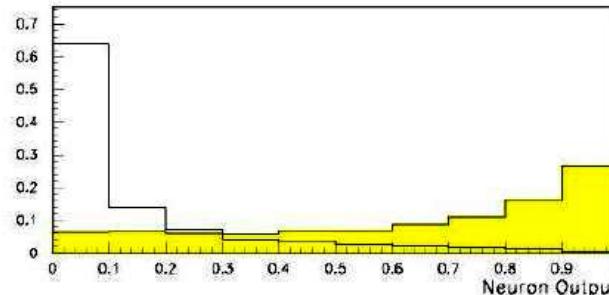
Signal:  $e^+e^- \rightarrow W^+W^-$  (often 4 well separated hadron jets)

Background:  $e^+e^- \rightarrow q\bar{q}gg$  (4 less well separated hadron jets)



← input variables based on jet structure, event shape, ...  
none by itself gives much separation.

Neural network output:



(Garrido, Juste and Martinez, ALEPH 96-144)

# Probability Density Estimation (PDE)

Construct non-parametric estimators for the pdfs of the data  $\mathbf{x}$  for the two event classes,  $p(\mathbf{x}|H_0), p(\mathbf{x}|H_1)$  and use these to construct the likelihood ratio, which we use for the discriminant function:

$$y(\mathbf{x}) = \frac{\hat{p}(\mathbf{x}|H_1)}{\hat{p}(\mathbf{x}|H_0)}$$

$n$ -dimensional histogram is a brute force example of this; we will see a number of ways that are much better.

# Correlation vs. independence

In general a multivariate distribution  $p(\mathbf{x})$  does not factorize into a product of the marginal distributions for the individual variables:

$$p(\vec{\mathbf{x}}) = \prod_{i=1}^n p_i(x_i)$$

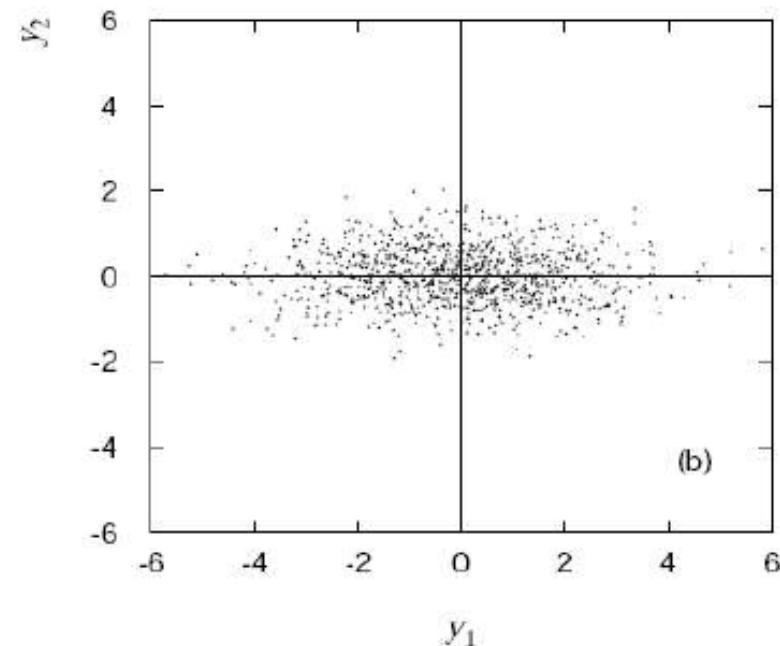
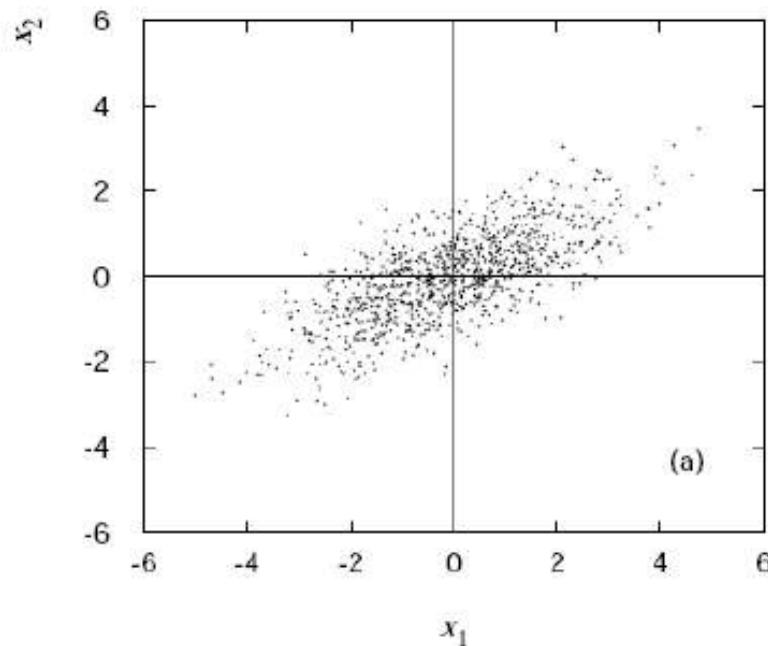
← holds only if the components of  $\mathbf{x}$  are independent

Most importantly, the components of  $\mathbf{x}$  will generally have nonzero covariances (i.e. they are correlated):

$$V_{ij} = \text{cov}[x_i, x_j] = E[x_i x_j] - E[x_i]E[x_j] \neq 0$$

# Decorrelation of input variables

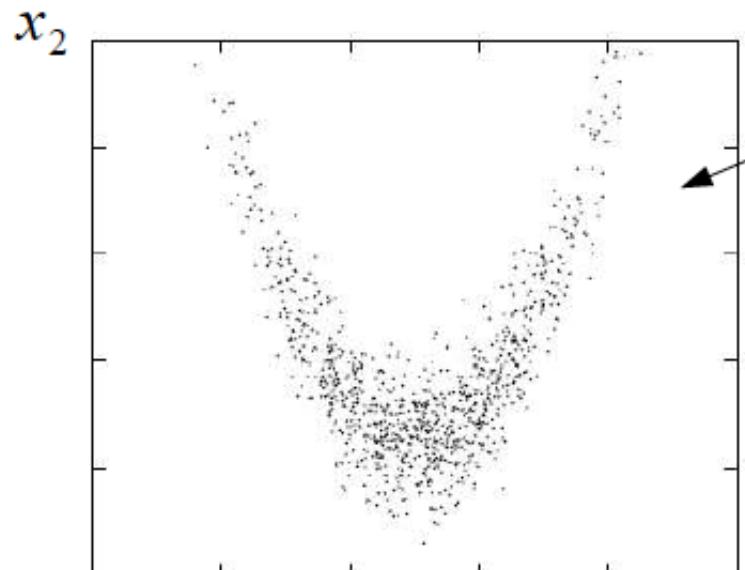
But we can define a set of uncorrelated input variables by a linear transformation, i.e., find the matrix  $A$  such that for  $\vec{y} = A \vec{x}$  the covariances  $\text{cov}[y_i, y_j] = 0$ :



For the following suppose that the variables are “decorrelated” in this way for each of  $p(\mathbf{x}|H_0)$  and  $p(\mathbf{x}|H_1)$  separately (since in general their correlations are different).

# Decorrelation is not enough

But even with zero correlation, a multivariate pdf  $p(\mathbf{x})$  will in general have nonlinearities and thus the decorrelated variables are still not independent.



pdf with zero covariance but  
components still not  
independent, since clearly

$$p(x_2|x_1) \equiv \frac{p(x_1, x_2)}{p_1(x_1)} \neq p_2(x_2)$$

and therefore

$$p(x_1, x_2) \neq p_1(x_1) p_2(x_2)$$

# Naive Bayes method

First decorrelate  $\mathbf{x}$ , i.e., find  $\mathbf{y} = A\mathbf{x}$ , with  $\text{cov}[y_i, y_j] = \text{V}[y_i] \delta_{ij}$ .

Pdfs of  $\mathbf{x}$  and  $\mathbf{y}$  are then related by

$$f(\mathbf{x}) = |J|g(\mathbf{y}(\mathbf{x})) \quad \text{where} \quad J = \det(A)$$

If nonlinear features of  $g(\mathbf{y})$  not too important, estimate using product of marginal pdfs:

$$\hat{f}(\mathbf{x}) = |J| \prod_{i=1}^n g_i(y_i(\mathbf{x})) = |\det(A)| \prod_{i=1}^n g_i((A\mathbf{x})_i)$$

Do separately for the two hypotheses s and b (separate matrices  $A_s$  and  $A_b$  and marginal pdfs  $g_{s,i}, g_{b,i}$ ). Then define test statistic as

$$y(\mathbf{x}) = \frac{\hat{f}_s(\mathbf{x})}{\hat{f}_b(\mathbf{x})}$$

Called Naive Bayes classifier. Reduces problem of estimating an  $n$ -dimensional pdf to finding  $n$  one-dimensional marginal pdfs.

# Kernel-based PDE (KDE)

Consider  $d$  dimensions,  $N$  training events,  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , estimate  $f(\mathbf{x})$  with

$$\hat{f}(\vec{x}) = \frac{1}{Nh^d} \sum_{i=1}^N K\left(\frac{\vec{x} - \vec{x}_i}{h}\right)$$

Annotations for the equation:

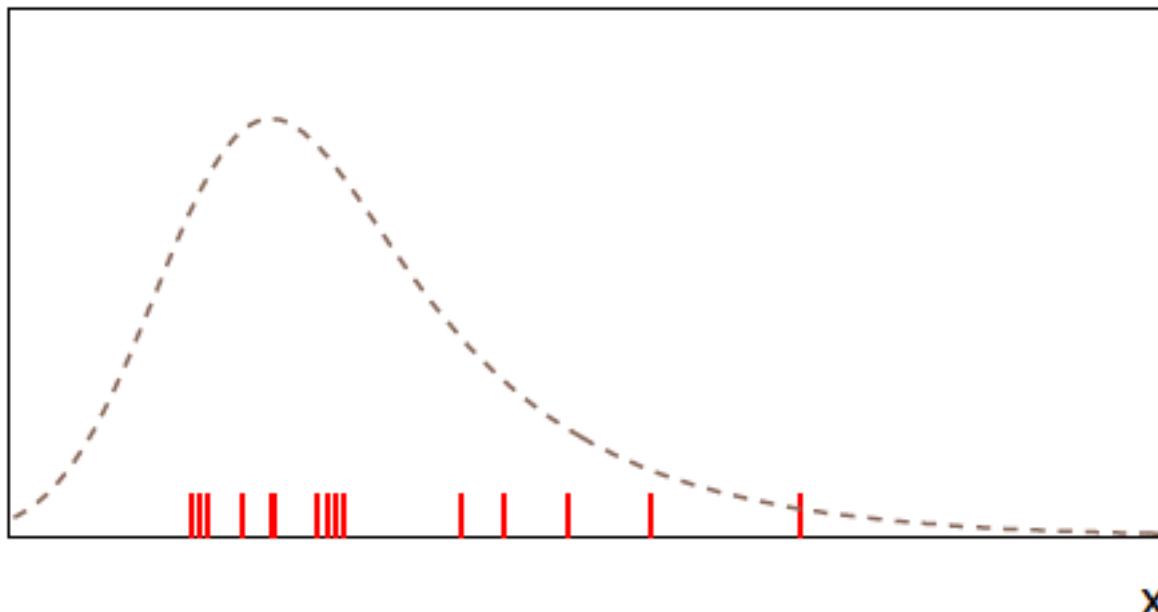
- $\mathbf{x}$  where we want to know pdf
- $\mathbf{x}$  of  $i^{\text{th}}$  training event
- bandwidth (smoothing parameter)
- kernel

Use e.g. Gaussian kernel:

$$K(\vec{x}) = \frac{1}{(2\pi)^{d/2}} e^{-|\vec{x}|^2/2}$$

# Gaussian KDE in 1-dimension

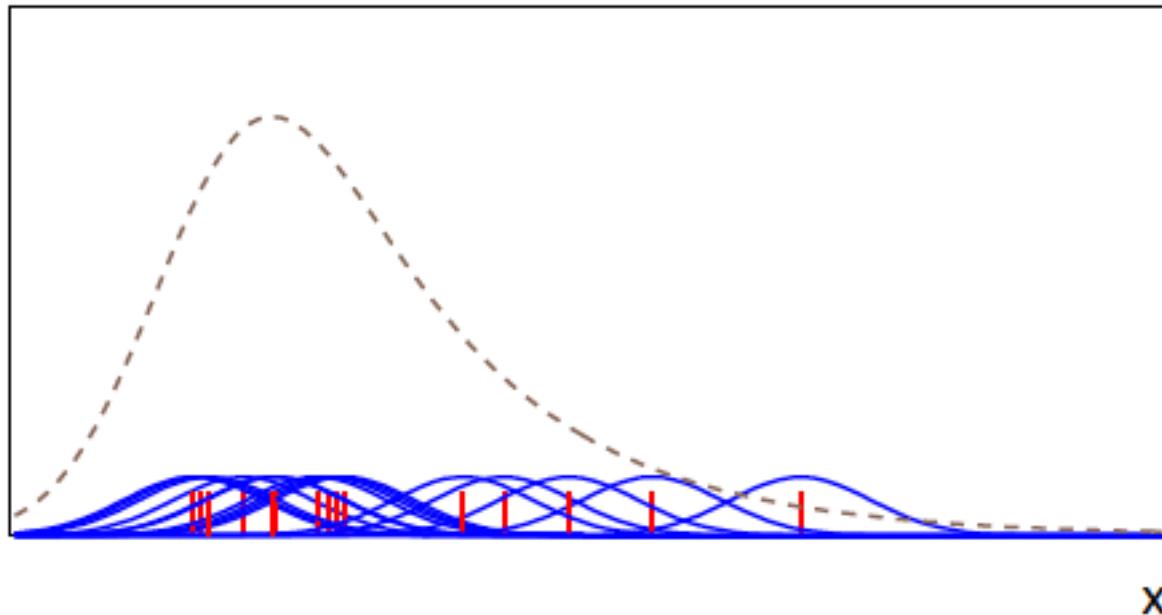
Suppose the pdf (dashed line) below is not known in closed form, but we can generate events that follow it (the red tick marks):



Goal is to find an approximation to the pdf using the generated data values.

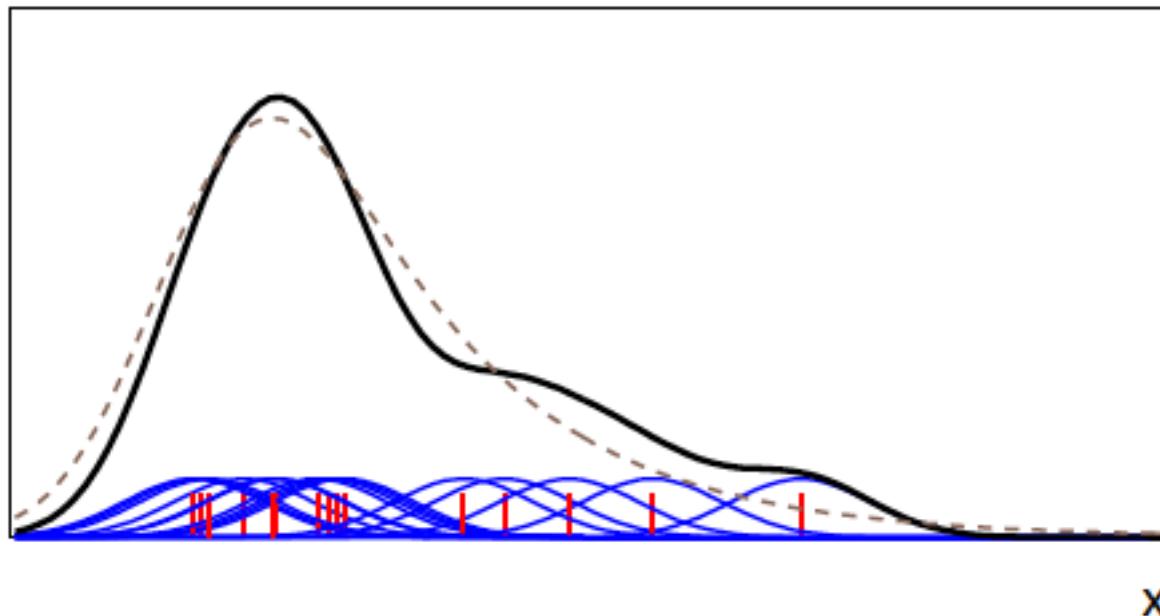
# Gaussian KDE in 1-dimension (cont.)

Place a kernel pdf (here a Gaussian) centred around each generated event weighted by  $1/N_{\text{event}}$ :



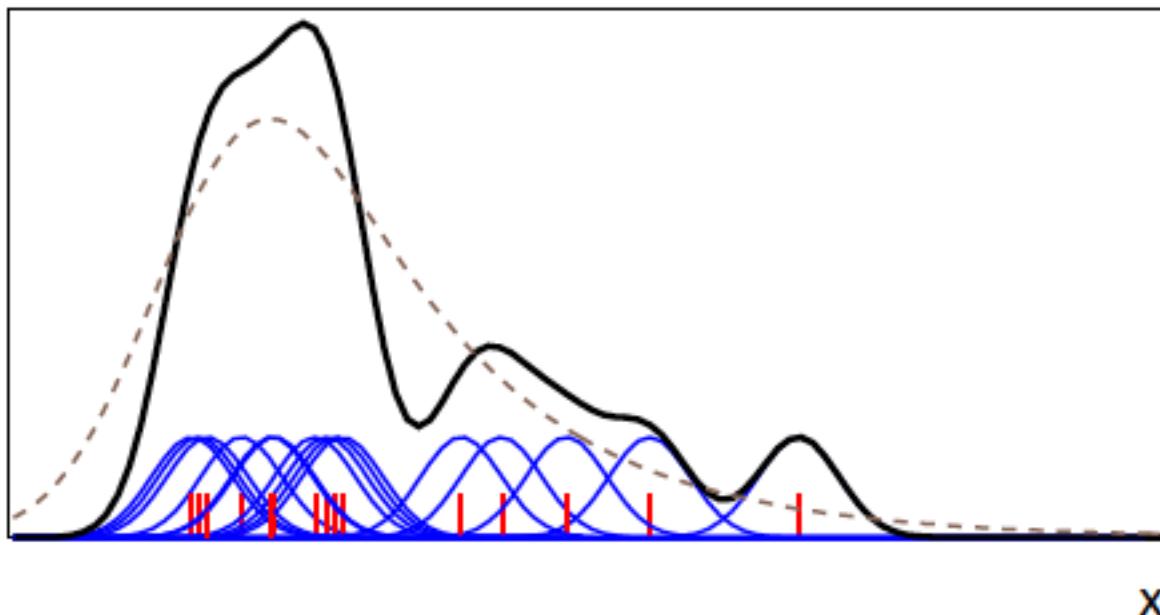
# Gaussian KDE in 1-dimension (cont.)

The KDE estimate the pdf is given by the sum of all of the Gaussians:



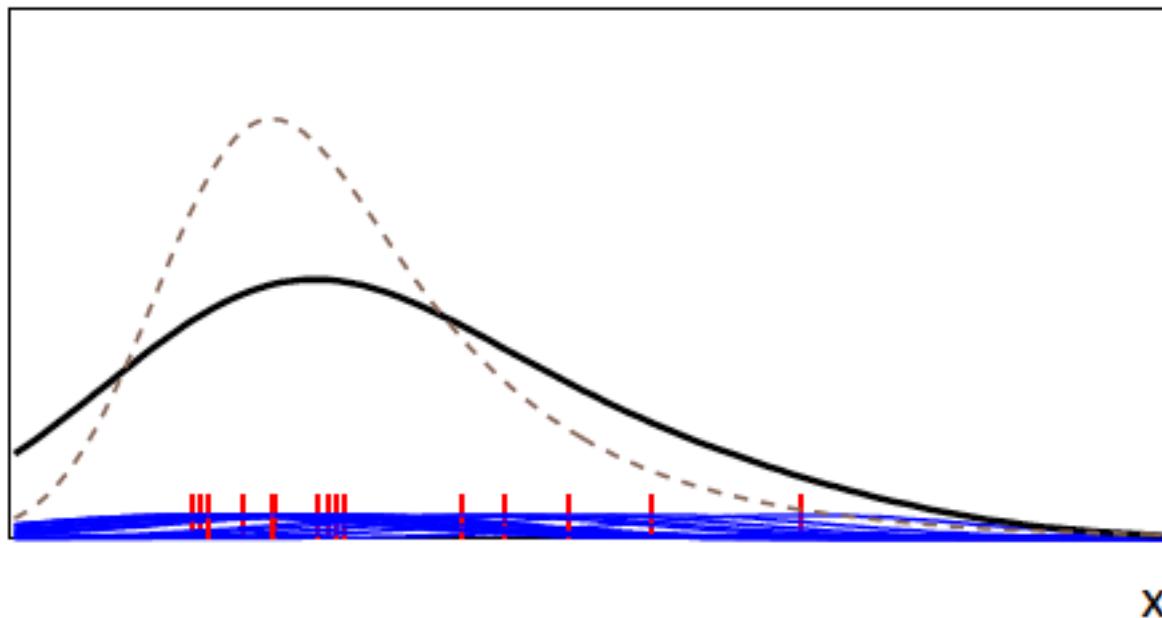
# Choice of kernel width

The width  $h$  of the Gaussians is analogous to the bin width of a histogram. If it is too small, the estimator has noise:



# Choice of kernel width (cont.)

If width of Gaussian kernels too large, structure is washed out:



# KDE discussion

Various strategies can be applied to choose width  $h$  of kernel based trade-off between bias and variance (noise).

Adaptive KDE allows width of kernel to vary, e.g., wide where target pdf is low (few events); narrow where pdf is high.

Advantage of KDE: no training!

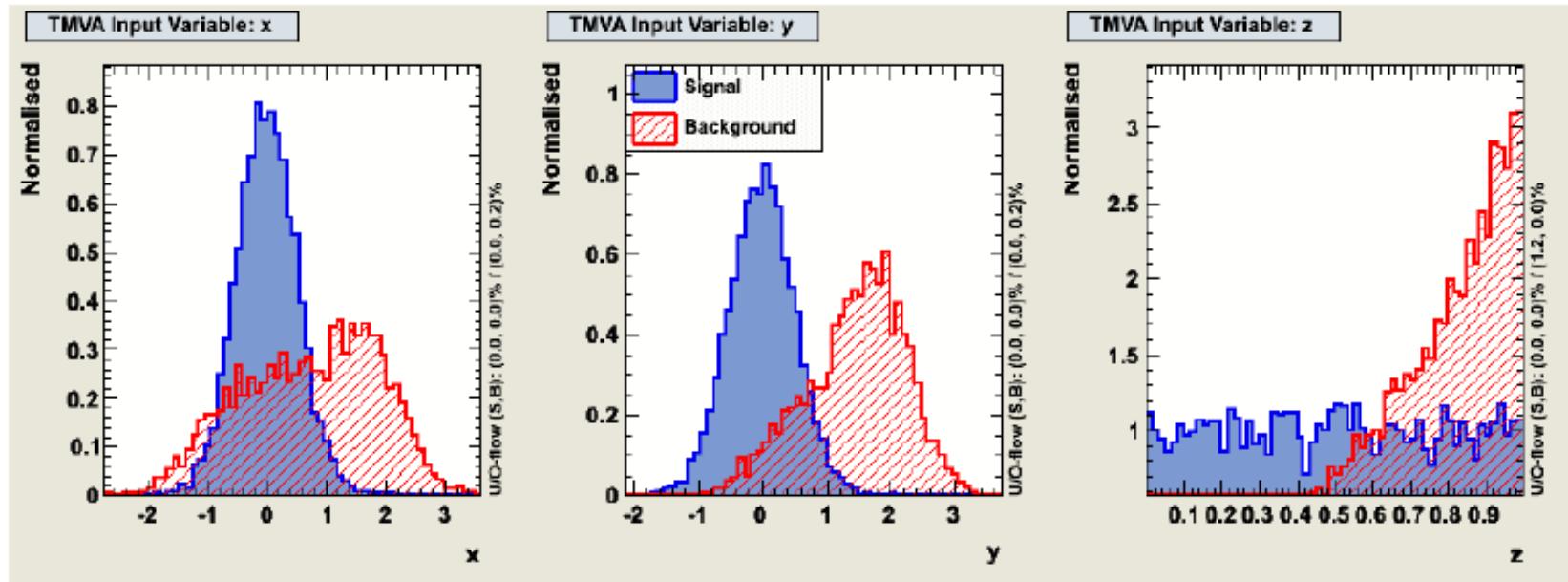
Disadvantage of KDE: to evaluate we need to sum  $N_{\text{event}}$  terms, so if we have many events this can be slow.

Special treatment required if kernel extends beyond range where pdf defined. Can e.g., renormalize the kernels to unity inside the allowed range; alternatively “mirror” the events about the boundary (contribution from the mirrored events exactly compensates the amount lost outside the boundary).

Software in ROOT: RooKeysPdf (K. Cranmer, CPC 136:198,2001)

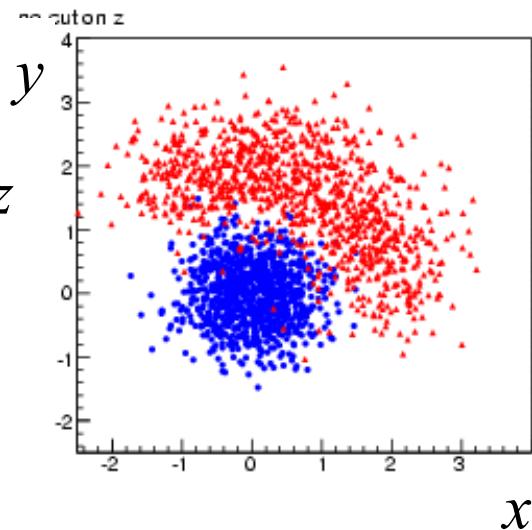
# Test example with TMVA

Each event characterized by 3 variables,  $x, y, z$ :

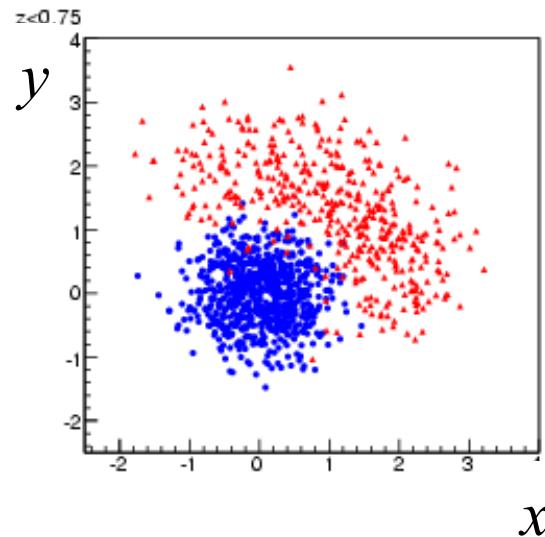


# Test example ( $x, y, z$ )

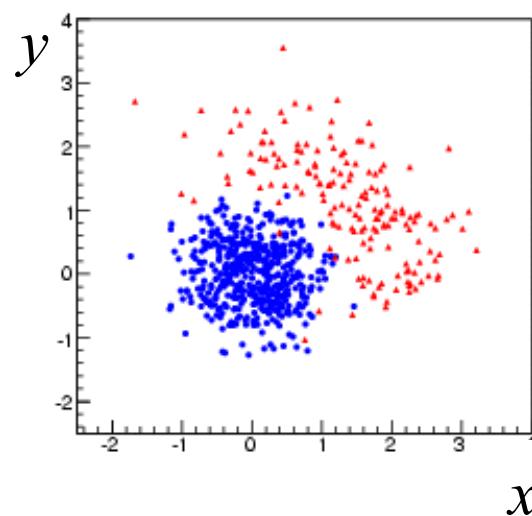
no cut on  $z$



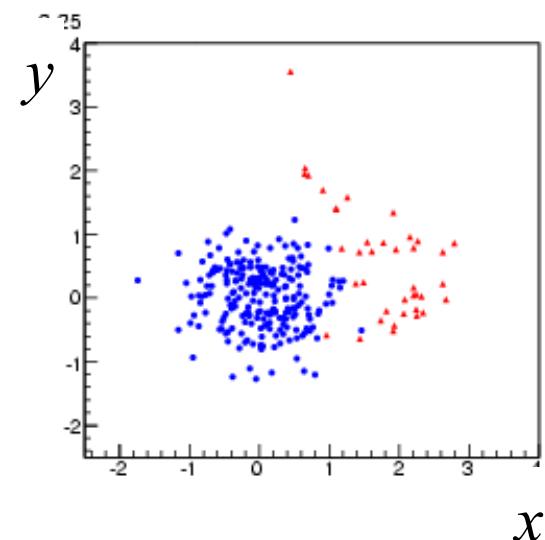
$z < 0.75$



$z < 0.5$



$z < 0.25$



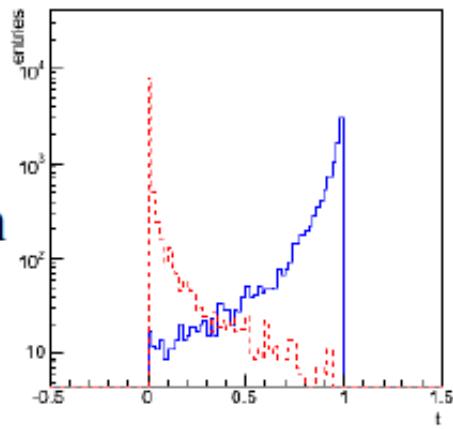
# Test example results

Fisher  
discriminant

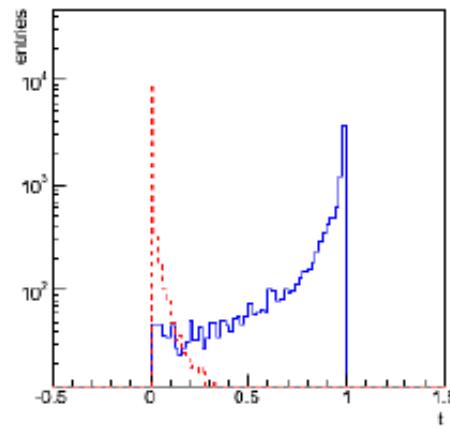
Multilayer  
perceptron

Find these on next homework assignment.

Naive Bayes,  
no decorrelation

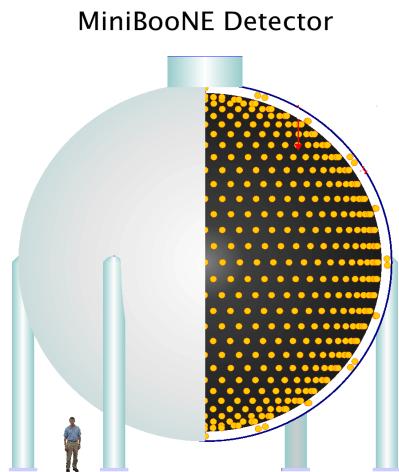


Naive Bayes with  
decorrelation

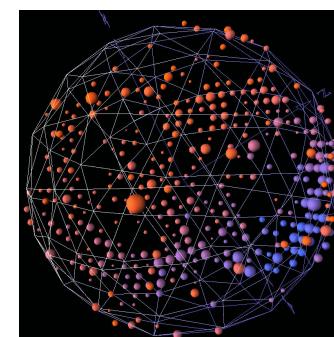
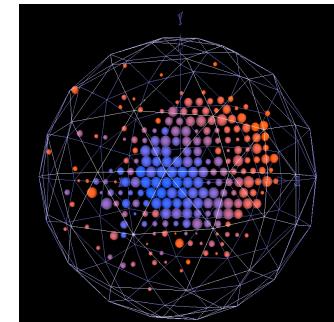
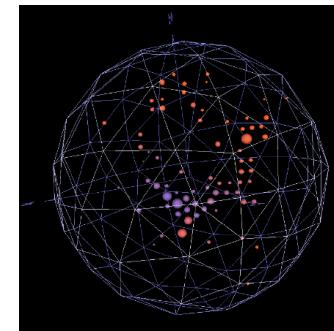
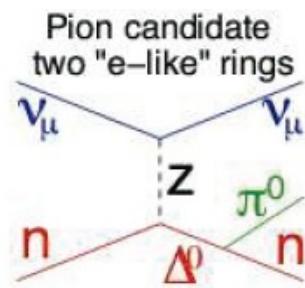
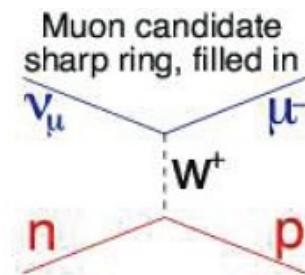
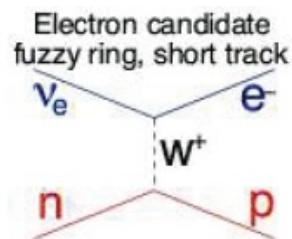


# Particle i.d. in MiniBooNE

Detector is a 12-m diameter tank of mineral oil exposed to a beam of neutrinos and viewed by 1520 photomultiplier tubes:



Search for  $\nu_\mu$  to  $\nu_e$  oscillations required particle i.d. using information from the PMTs.



H.J. Yang, MiniBooNE PID, DNP06

# Decision trees

Out of all the input variables, find the one for which with a single cut gives best improvement in signal purity:

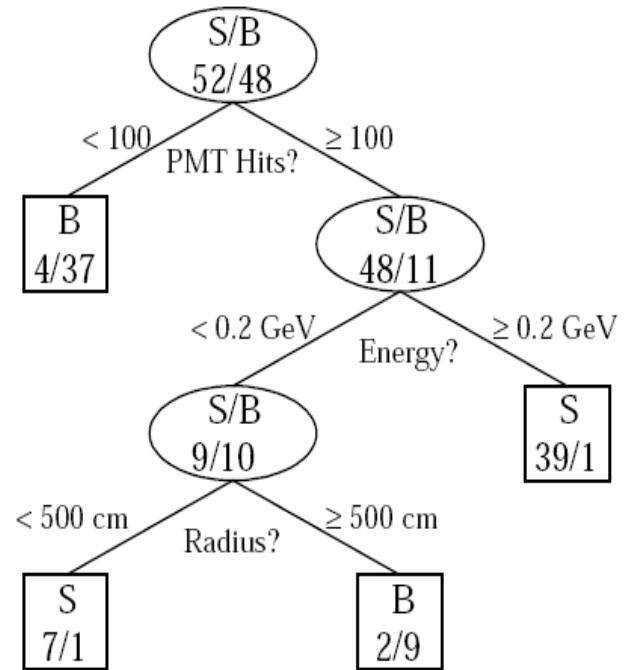
$$P = \frac{\sum_{\text{signal}} w_i}{\sum_{\text{signal}} w_i + \sum_{\text{background}} w_i}$$

where  $w_i$  is the weight of the  $i$ th event.

Resulting nodes classified as either signal/background.

Iterate until stop criterion reached based on e.g. purity or minimum number of events in a node.

The set of cuts defines the decision boundary.



Example by MiniBooNE experiment,  
B. Roe et al., NIM 543 (2005) 577

# Finding the best single cut

The level of separation within a node can, e.g., be quantified by the *Gini coefficient*, calculated from the (s or b) purity as:

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

For a cut that splits a set of events  $a$  into subsets  $b$  and  $c$ , one can quantify the improvement in separation by the change in weighted Gini coefficients:

$$\Delta = W_a G_a - W_b G_b - W_c G_c \quad \text{where, e.g.,} \quad W_a = \sum_{i \in a} w_i$$

Choose e.g. the cut to the maximize  $\Delta$ ; a variant of this scheme can use instead of Gini e.g. the misclassification rate:

$$\varepsilon = 1 - \max(p, 1 - p)$$

# Decision trees (2)

The terminal nodes (**leaves**) are classified a signal or background depending on majority vote (or e.g. signal fraction greater than a specified threshold).

This classifies every point in input-variable space as either signal or background, a **decision tree classifier**, with discriminant function

$$f(\mathbf{x}) = 1 \text{ if } \mathbf{x} \text{ in signal region, } -1 \text{ otherwise}$$

Decision trees tend to be very sensitive to statistical fluctuations in the training sample.

Methods such as **boosting** can be used to stabilize the tree.

# Boosting

Boosting is a general method of creating a set of classifiers which can be combined to achieve a new classifier that is more stable and has a smaller error than any individual one.

Often applied to decision trees but, can be applied to any classifier.

Suppose we have a training sample  $T$  consisting of  $N$  events with

- $\mathbf{x}_1, \dots, \mathbf{x}_N$  event data vectors (each  $\mathbf{x}$  multivariate)
- $y_1, \dots, y_N$  true class labels, +1 for signal, -1 for background
- $w_1, \dots, w_N$  event weights

Now define a rule to create from this an ensemble of training samples  $T_1, T_2, \dots$ , derive a classifier from each and average them.

Trick is to create modifications in the training sample that give classifiers with smaller error rates than those of the preceding ones.

A successful example is AdaBoost (Freund and Schapire, 1997).

# AdaBoost

First initialize the training sample  $T_1$  using the original

$\mathbf{x}_1, \dots, \mathbf{x}_N$  event data vectors

$y_1, \dots, y_N$  true class labels (+1 or -1)

$w_1^{(1)}, \dots, w_N^{(1)}$  event weights

with the weights equal and normalized such that

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

Then train the classifier  $f_1(\mathbf{x})$  (e.g. a decision tree) with a method that incorporates the event weights. For an event with data  $\mathbf{x}_i$ ,

$f_1(\mathbf{x}_i) > 0$  classify as signal

$f_1(\mathbf{x}_i) < 0$  classify as background

# Updating the event weights

Define the training sample for step  $k+1$  from that of  $k$  by updating the event weights according to

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

$i$  = event index

$k$  = training sample index

where  $Z_k$  is a normalization factor defined such that the sum of the weights over all events is equal to one.

Therefore event weight for event  $i$  is increased in the  $k+1$  training sample if it was classified incorrectly in sample  $k$ .

Idea is that next time around the classifier should pay more attention to this event and try to get it right.

# Error rate of the $k$ th classifier

At each step the classifiers  $f_k(\mathbf{x})$  are defined so as to minimize the error rate  $\varepsilon_k$ ,

$$\varepsilon_k = \sum_{i=1}^N w_i^{(k)} I(y_i f_k(\mathbf{x}_i) \leq 0)$$

where  $I(X) = 1$  if  $X$  is true and is zero otherwise.

# Assigning the classifier score

Assign a score to the  $k$ th classifier based on its error rate,

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

If we define the final classifier as  $f(\mathbf{x}) = \sum_{k=1}^K \alpha_k f_k(\mathbf{x}, T_k)$

then one can show that its error rate on the training data satisfies the bound

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

# AdaBoost error rate

So providing each classifier in the ensemble has  $\varepsilon_k < \frac{1}{2}$ , i.e., better than random guessing, then the error rate for the final classifier on the training data (not on unseen data) drops to zero.

That is, for sufficiently large  $K$  the training data will be over fitted.

The error rate on a validation sample would reach some minimum after a certain number of steps and then could rise.

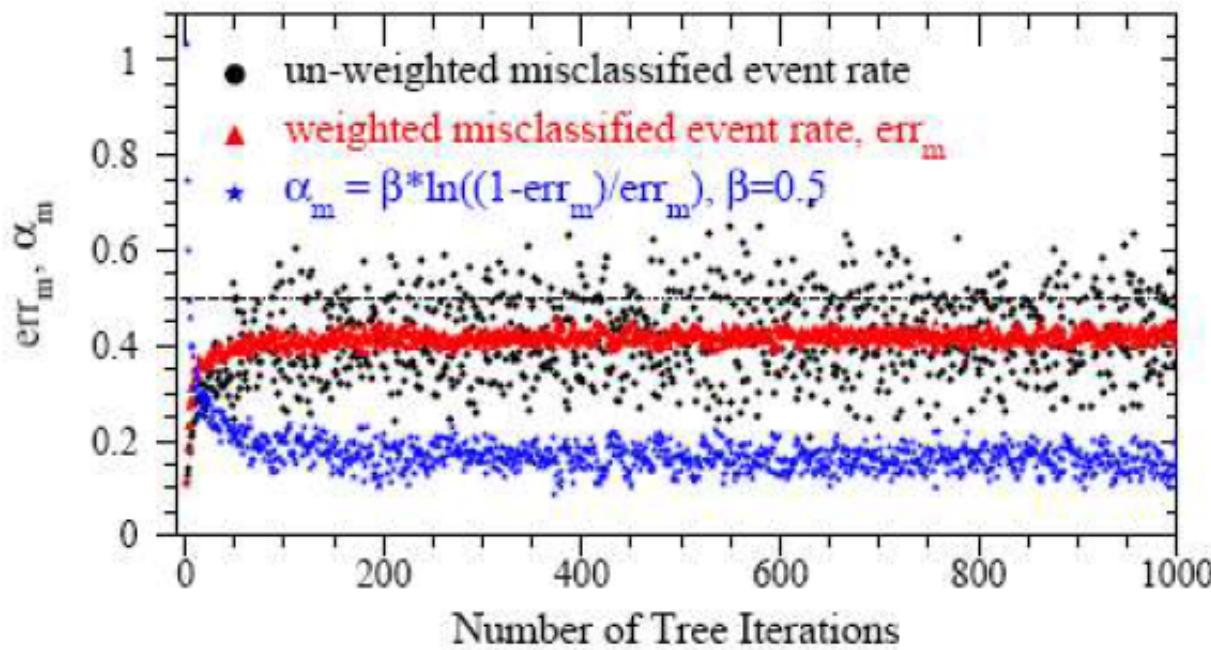
So the procedure is to monitor the error rate of the combined classifier at each step with a validation sample and to stop before it starts to rise.

Although in principle AdaBoost must overfit, in practice following this procedure overtraining is not a big problem.

# BDT example from MiniBooNE

~200 input variables for each event ( $\nu$  interaction producing e,  $\mu$  or  $\pi$ ).

Each individual tree is relatively weak, with a misclassification error rate  $\sim 0.4 - 0.45$

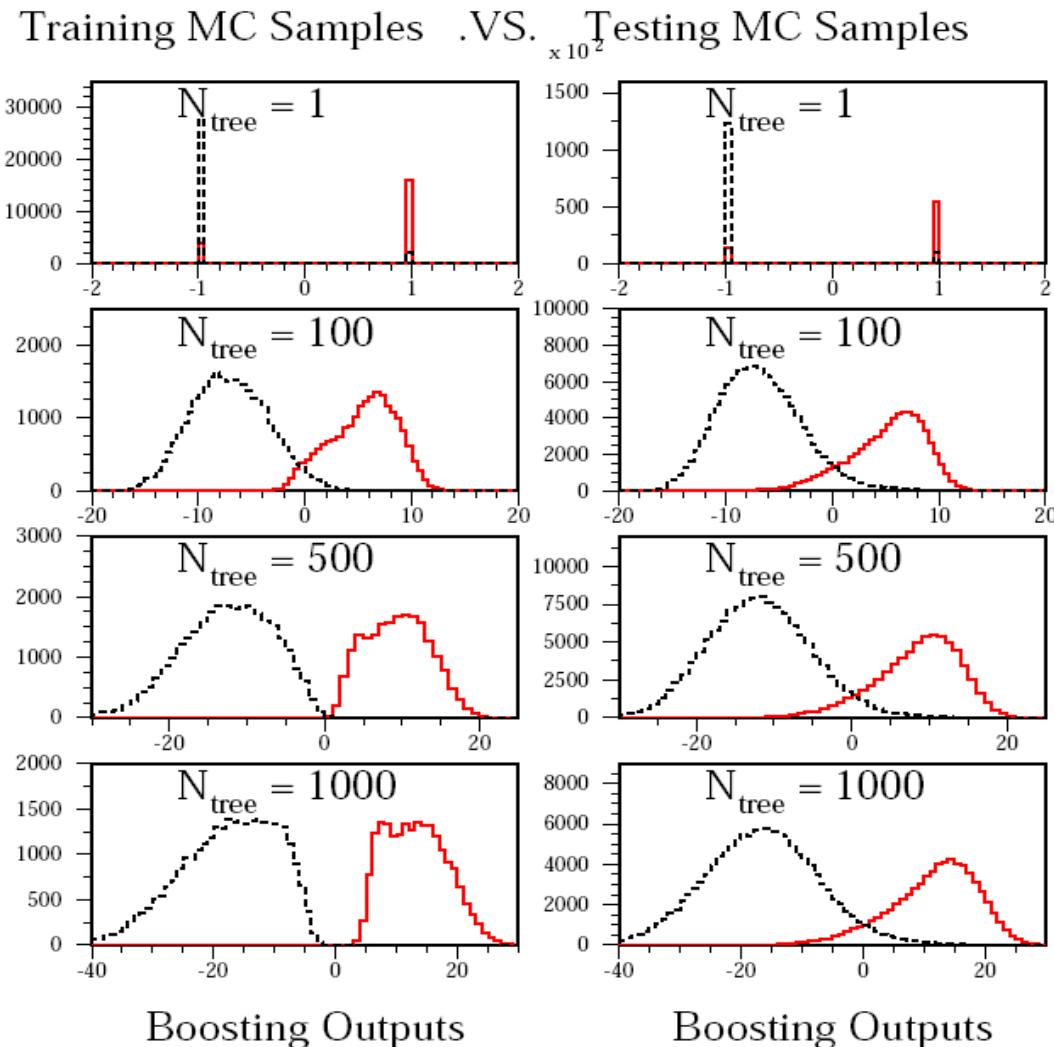


B. Roe et al., NIM 543 (2005) 577

# Monitoring overtraining

From MiniBooNE example:

Performance stable after a few hundred trees.



# A simple example (2D)

Consider two variables,  $x_1$  and  $x_2$ , and suppose we have formulas for the joint pdfs for both signal (s) and background (b) events (in real problems the formulas are usually not available).

$f(x_1|x_2) \sim \text{Gaussian}$ , different means for s/b,

Gaussians have same  $\sigma$ , which depends on  $x_2$ ,

$f(x_2) \sim \text{exponential}$ , same for both s and b,

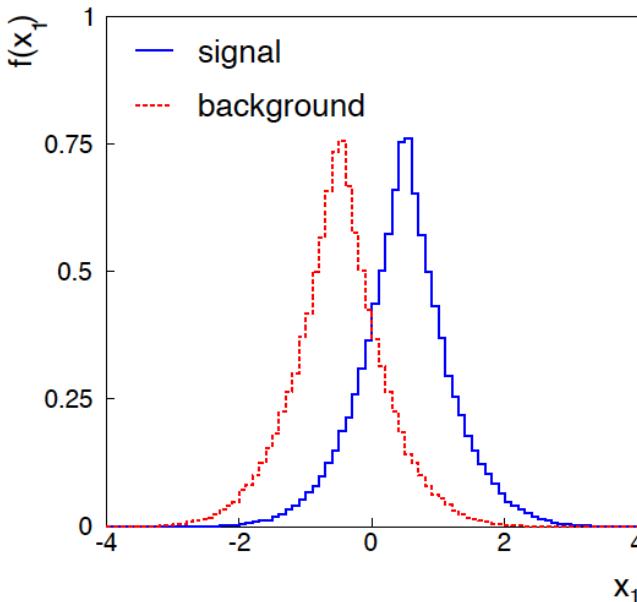
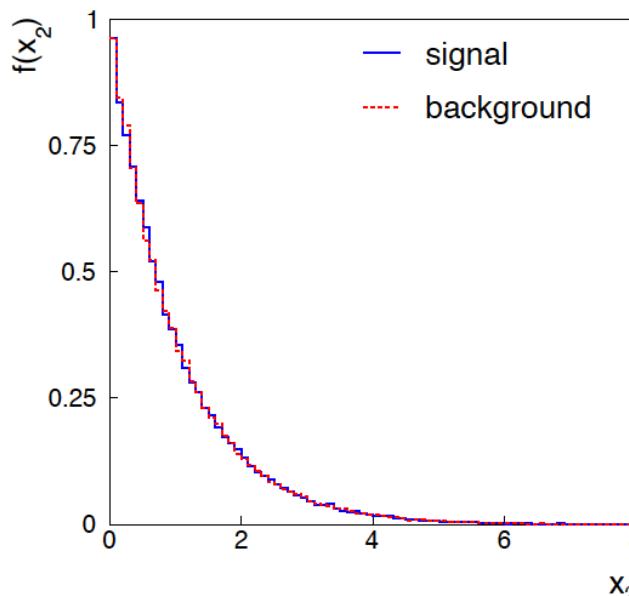
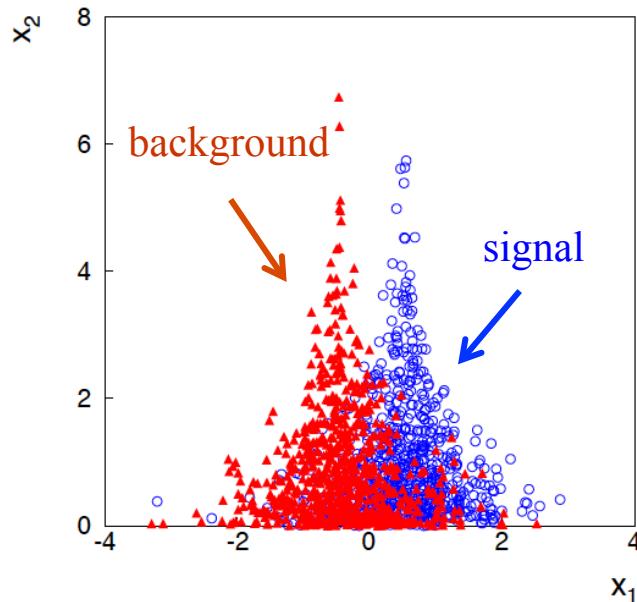
$f(x_1, x_2) = f(x_1|x_2) f(x_2)$ :

$$f(x_1, x_2|s) = \frac{1}{\sqrt{2\pi}\sigma(x_2)} e^{-(x_1 - \mu_s)^2/2\sigma^2(x_2)} \frac{1}{\lambda} e^{-x_2/\lambda}$$

$$f(x_1, x_2|b) = \frac{1}{\sqrt{2\pi}\sigma(x_2)} e^{-(x_1 - \mu_b)^2/2\sigma^2(x_2)} \frac{1}{\lambda} e^{-x_2/\lambda}$$

$$\sigma(x_2) = \sigma_0 e^{-x_2/\xi}$$

# Joint and marginal distributions of $x_1, x_2$



Distribution  $f(x_2)$  same for s, b.

So does  $x_2$  help discriminate  
between the two event types?

# Likelihood ratio for 2D example

Neyman-Pearson lemma says best critical region is determined by the likelihood ratio:

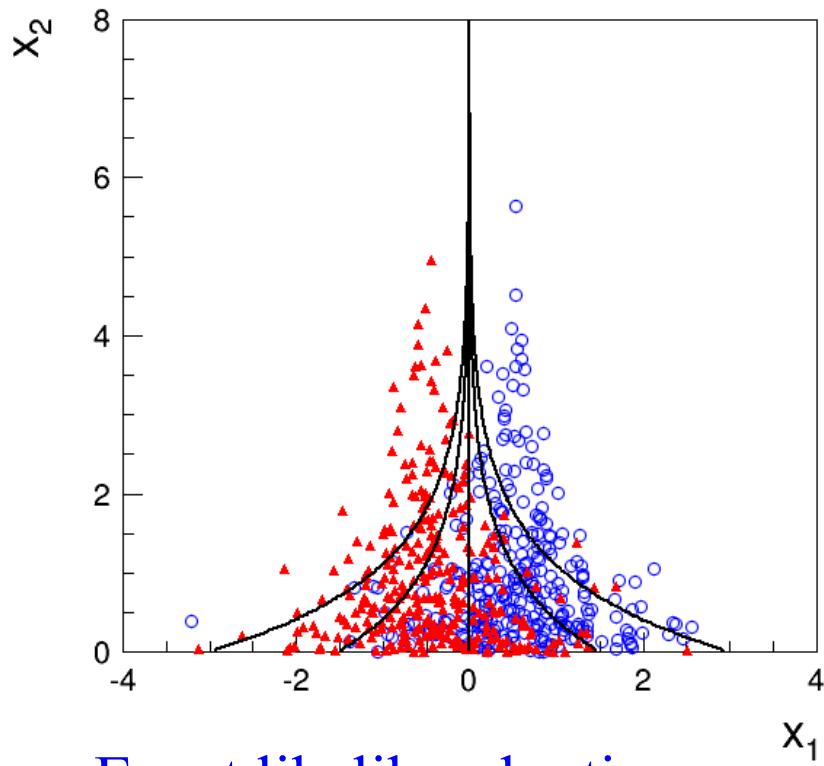
$$t(x_1, x_2) = \frac{f(x_1, x_2 | s)}{f(x_1, x_2 | b)}$$

Equivalently we can use any monotonic function of this as a test statistic, e.g.,

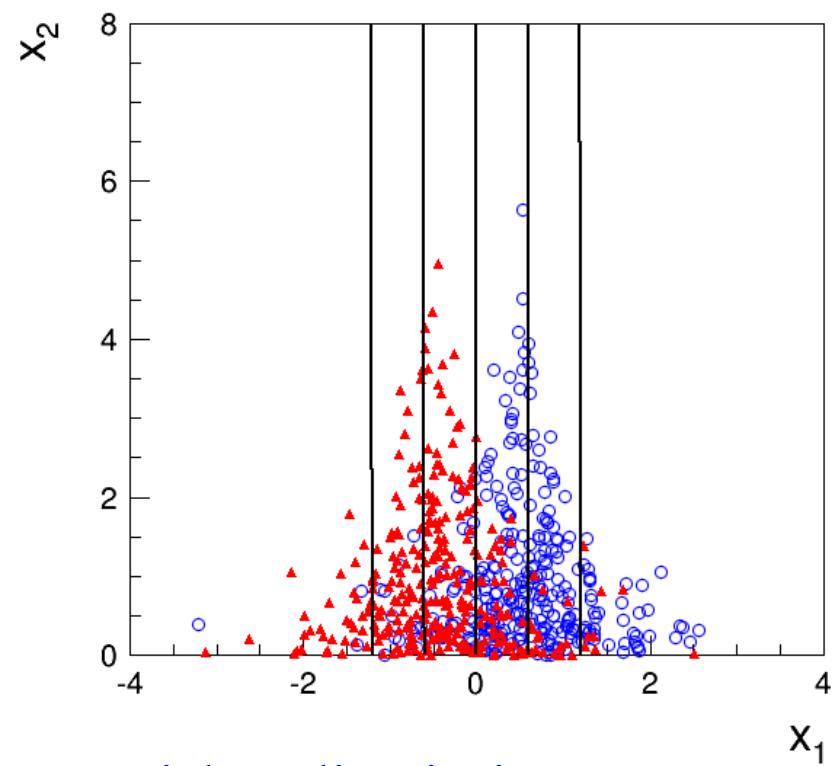
$$\ln t = \frac{\frac{1}{2}(\mu_b^2 - \mu_s^2) + (\mu_s - \mu_b)x_1}{\sigma_0^2 e^{-2x_2/\xi}}$$

Boundary of optimal critical region will be curve of constant  $\ln t$ , and this depends on  $x_2$ !

# Contours of constant MVA output

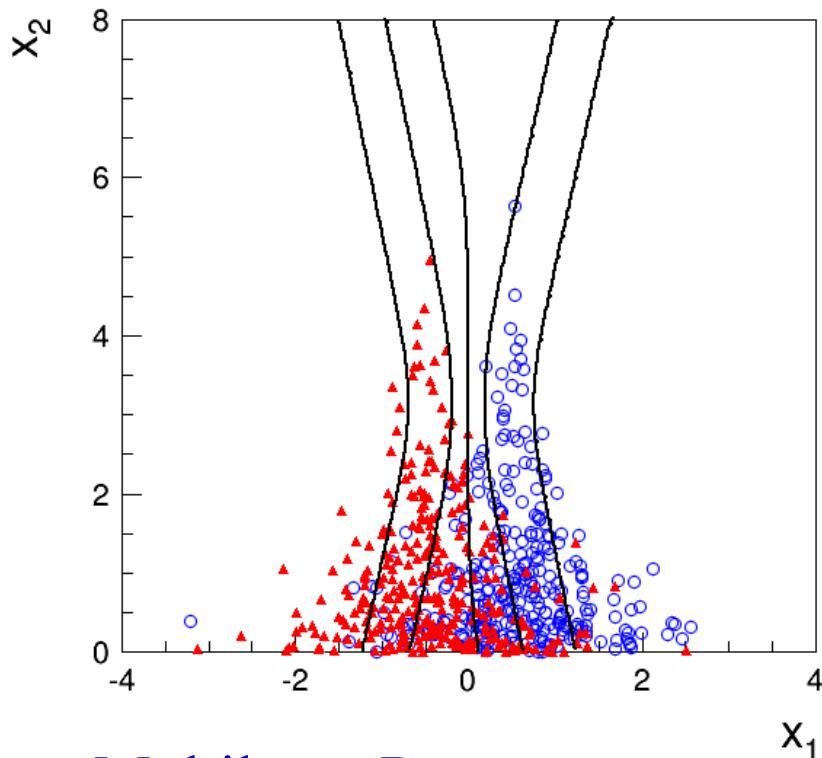


Exact likelihood ratio

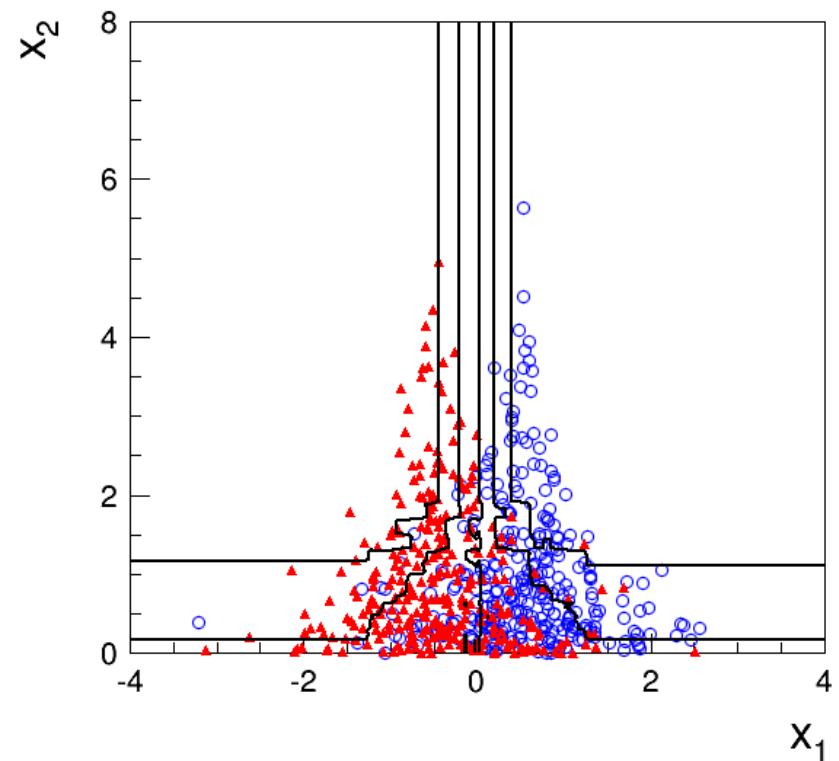


Fisher discriminant

# Contours of constant MVA output



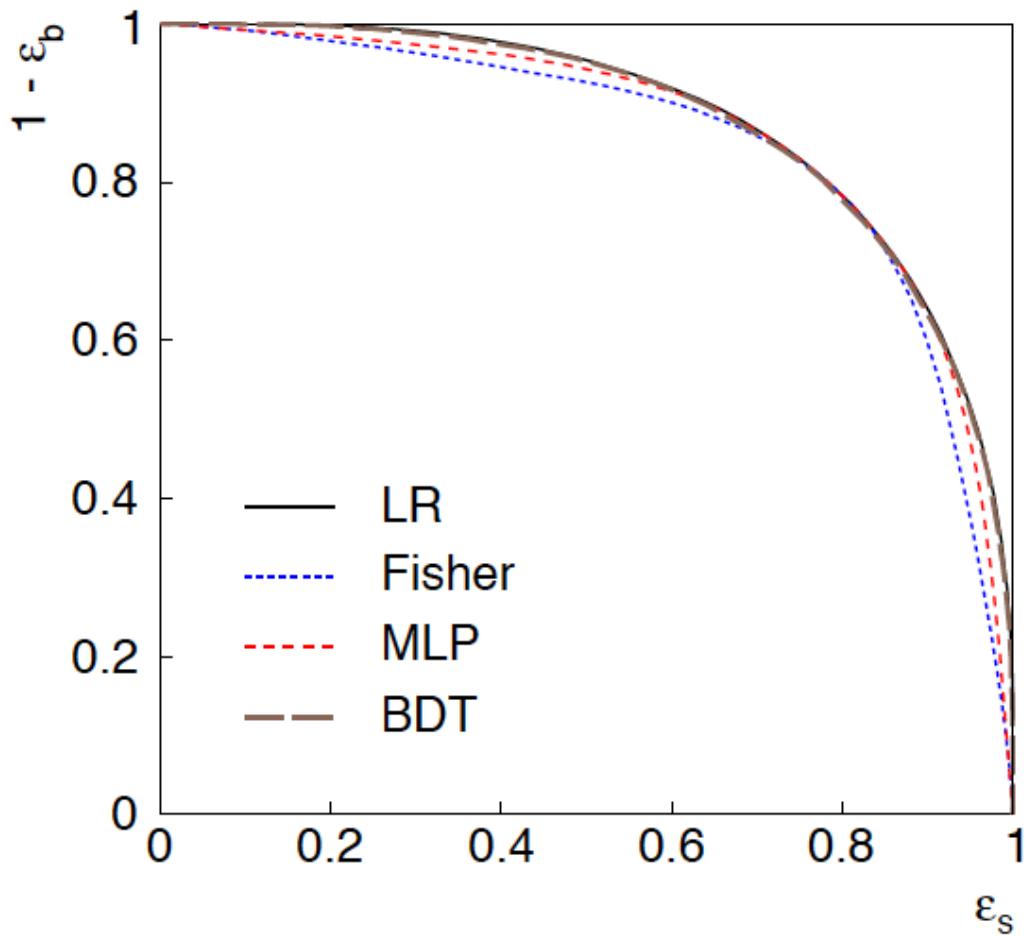
Multilayer Perceptron  
1 hidden layer with 2 nodes



Boosted Decision Tree  
200 iterations (AdaBoost)

Training samples:  $10^5$  signal and  $10^5$  background events

# ROC curve



ROC = “receiver operating characteristic” (term from signal processing).

Shows (usually) background rejection ( $1 - \varepsilon_b$ ) versus signal efficiency  $\varepsilon_s$ .

Higher curve is better; usually analysis focused on a small part of the curve.

## 2D Example: discussion

Even though the distribution of  $x_2$  is same for signal and background,  $x_1$  and  $x_2$  are not independent, so using  $x_2$  as an input variable helps.

Here we can understand why: high values of  $x_2$  correspond to a smaller  $\sigma$  for the Gaussian of  $x_1$ . So high  $x_2$  means that the value of  $x_1$  was well measured.

If we don't consider  $x_2$ , then all of the  $x_1$  measurements are lumped together. Those with large  $\sigma$  (low  $x_2$ ) "pollute" the well measured events with low  $\sigma$  (high  $x_2$ ).

Often in HEP there may be variables that are characteristic of how well measured an event is (region of detector, number of pile-up vertices,...). Including these variables in a multivariate analysis preserves the information carried by the well-measured events, leading to improved performance.

# Summary on multivariate methods

Particle physics has used several multivariate methods for many years:

- linear (Fisher) discriminant

- neural networks

- naive Bayes

and has in recent years started to use a few more:

- boosted decision trees

- support vector machines

- kernel density estimation

- $k$ -nearest neighbour

The emphasis is often on controlling systematic uncertainties between the modeled training data and Nature to avoid false discovery.

Although many classifier outputs are "black boxes", a discovery at  $5\sigma$  significance with a sophisticated (opaque) method will win the competition if backed up by, say,  $4\sigma$  evidence from a cut-based method.