

Topics in Statistical Data Analysis for HEP

Lecture 2: Statistical Tests / Multivariate Methods



CERN – Latin-American School
on High Energy Physics

Natal, Brazil, 4 April 2011



Glen Cowan
Physics Department
Royal Holloway, University of London
g.cowan@rhul.ac.uk
www.pp.rhul.ac.uk/~cowan

Outline

Lecture 1: Introduction and basic formalism

Probability

Parameter estimation

Statistical tests



Lecture 2: Statistics for making a discovery

Multivariate methods

Discovery significance and sensitivity

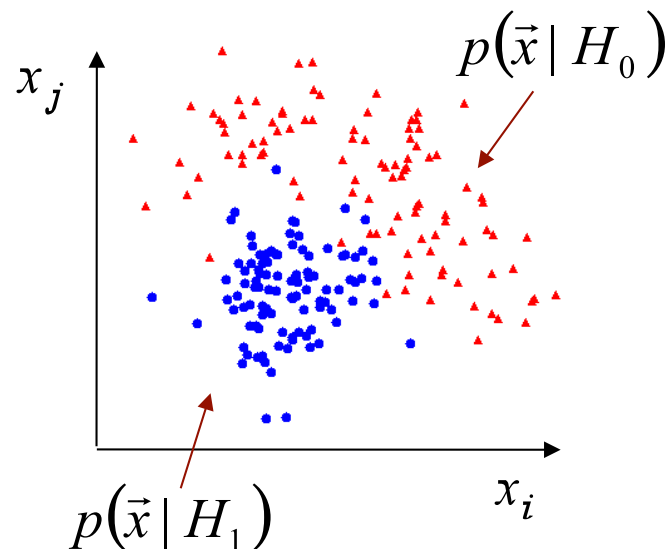
Systematic uncertainties

Event selection as a statistical test

For each event we measure a set of numbers: $\vec{x} = (x_1, \dots, x_n)$

$x_1 = \text{jet } p_T$
 $x_2 = \text{missing energy}$
 $x_3 = \text{particle i.d. measure, ...}$

\vec{x} follows some n -dimensional joint probability density, which depends on the type of event produced, i.e., was it $pp \rightarrow t\bar{t}$, $pp \rightarrow \tilde{g}\tilde{g}, \dots$

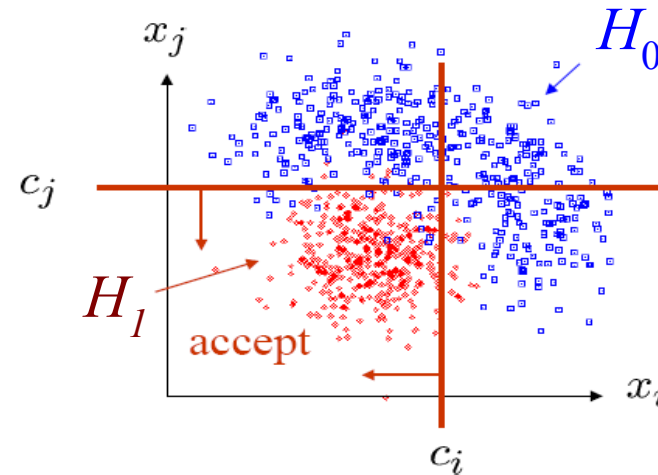


E.g. hypotheses H_0, H_1, \dots
Often simply “signal”,
“background”

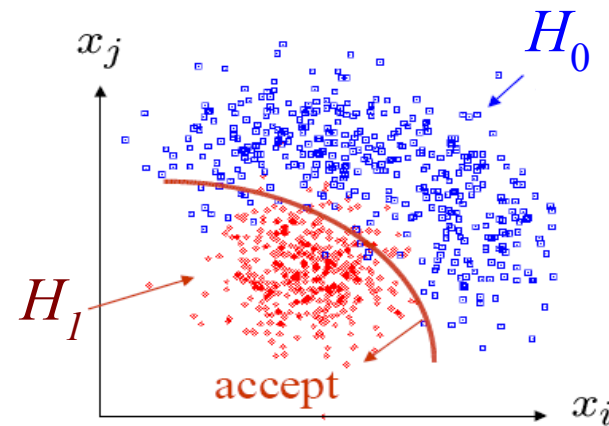
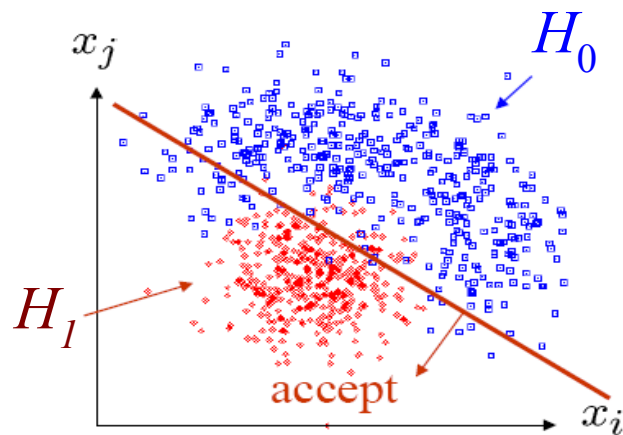
Finding an optimal decision boundary

In particle physics usually start by making simple “cuts”:

$$\begin{aligned}x_i &< c_i \\x_j &< c_j\end{aligned}$$



Maybe later try some other type of decision boundary:



Test statistics

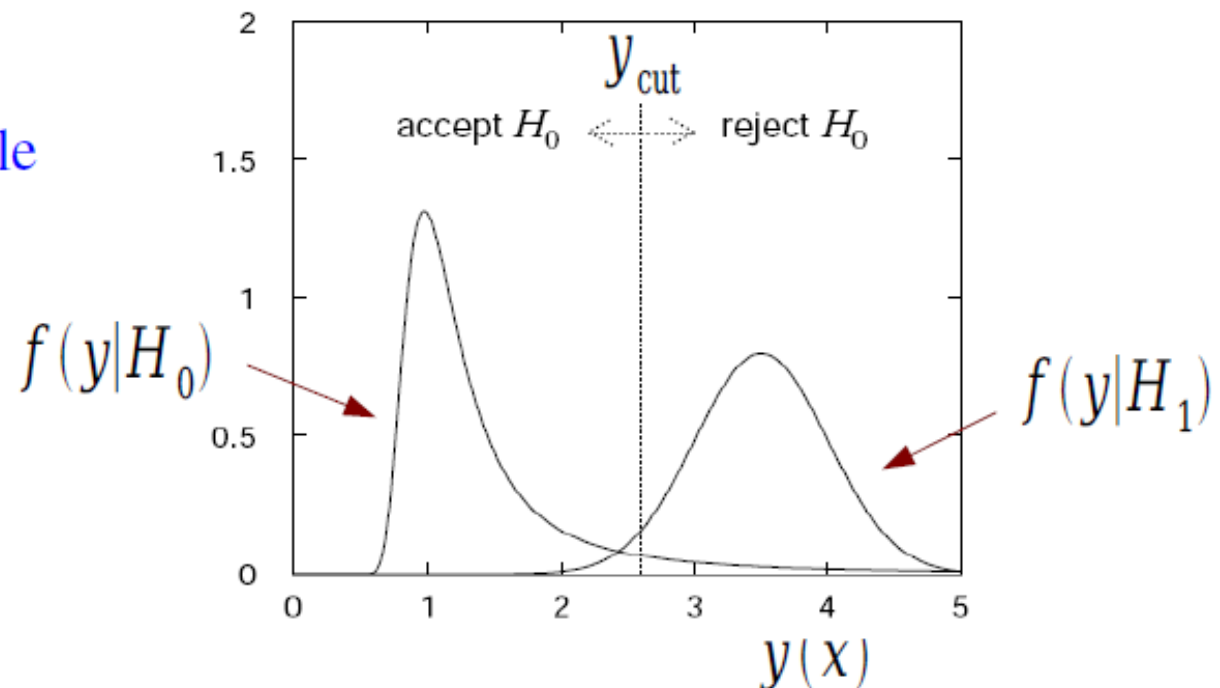
The decision boundary is a surface in the n -dimensional space of input variables, e.g., $y(\vec{x}) = \text{const.}$

We can treat the $y(x)$ as a scalar **test statistic** or discriminating function, and try to define this function so that its distribution has the maximum possible separation between the event types:

The decision boundary is now effectively a single cut on $y(x)$, dividing x -space into two regions:

R_0 (accept H_0)

R_1 (reject H_0)



Constructing a test statistic

The **Neyman-Pearson lemma** states: to obtain the highest background rejection for a given signal efficiency (highest power for a given significance level), choose the acceptance region for signal such that

$$\frac{p(\vec{x}|\mathbf{s})}{p(\vec{x}|\mathbf{b})} > c$$

where c is a constant that determines the signal efficiency.

Equivalently, the optimal discriminating function is given by the **likelihood ratio**:

$$y(\vec{x}) = \frac{p(\vec{x}|\mathbf{s})}{p(\vec{x}|\mathbf{b})}$$

N.B. any monotonic function of this is just as good.

Neyman-Pearson doesn't always help

The problem is that we usually don't have explicit formulae for the pdfs $p(\mathbf{x}|\mathbf{s})$, $p(\mathbf{x}|\mathbf{b})$, so for a given \mathbf{x} we can't evaluate the likelihood ratio.

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

$$\begin{array}{lll} \text{generate } \vec{x} \sim p(\vec{x}|\mathbf{s}) & \longrightarrow & \vec{x}_1, \dots, \vec{x}_{N_s} \\ \text{generate } \vec{x} \sim p(\vec{x}|\mathbf{b}) & \longrightarrow & \vec{x}_1, \dots, \vec{x}_{N_b} \end{array} \quad \begin{array}{l} \text{“training data”} \\ \text{events of known type} \end{array}$$

Naive try: enter each (s,b) event into an n -dimensional histogram, use e.g. M bins for each of the n dimensions, total of M^n cells.

n is potentially large \rightarrow prohibitively large number of cells to populate, can't generate enough training data.

Some “standard” multivariate methods

Place cuts on individual variables

Simple, intuitive, in general not optimal

Linear discriminant (e.g. Fisher)

Simple, optimal if the event types are Gaussian distributed with equal covariance, otherwise not optimal.

Probability Density Estimation based methods

Try to estimate $p(\mathbf{x}|s)$, $p(\mathbf{x}|b)$ then use $y(\vec{x}) = \hat{p}(\mathbf{x}|s) / \hat{p}(\mathbf{x}|b)$.

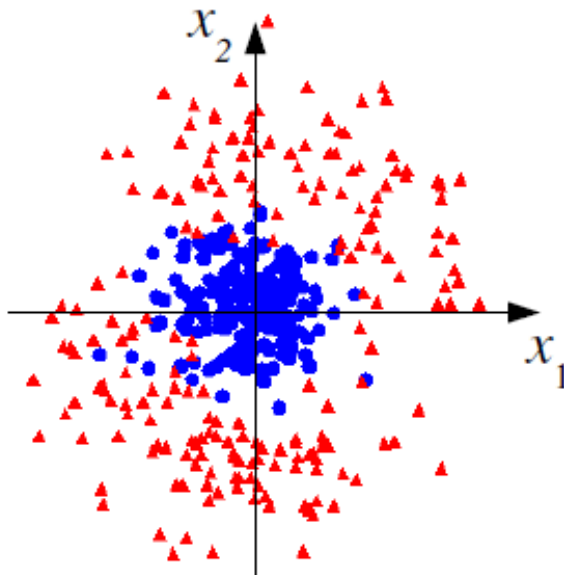
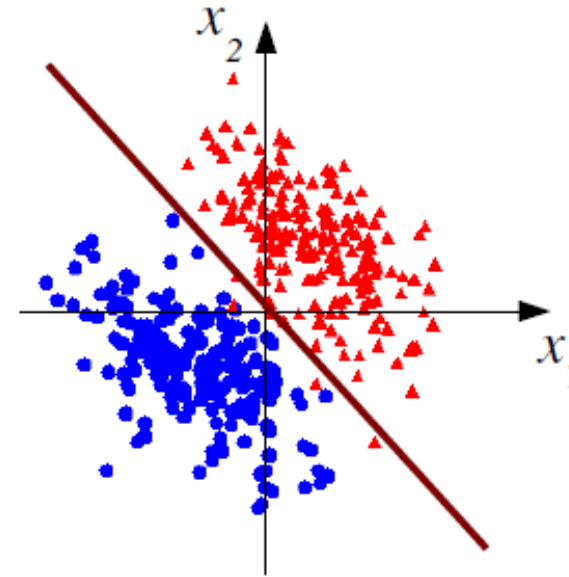
In principle best, difficult to estimate $p(\mathbf{x})$ for high dimension.

Neural networks

Can produce arbitrary decision boundary (in principle optimal), but can be difficult to train, result non-intuitive.

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.

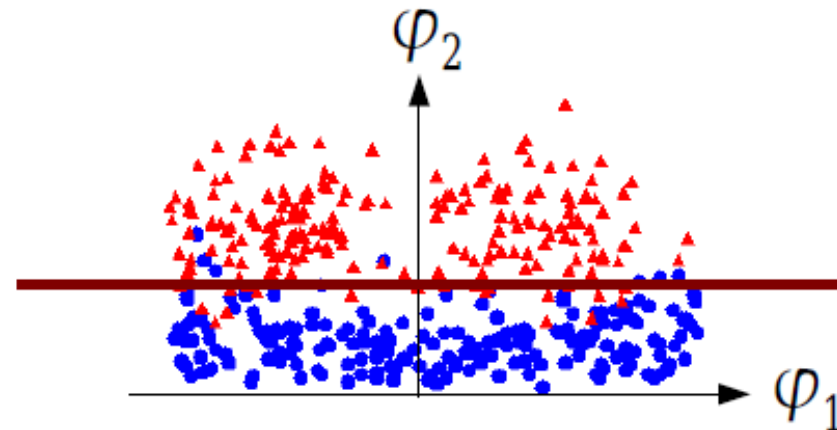
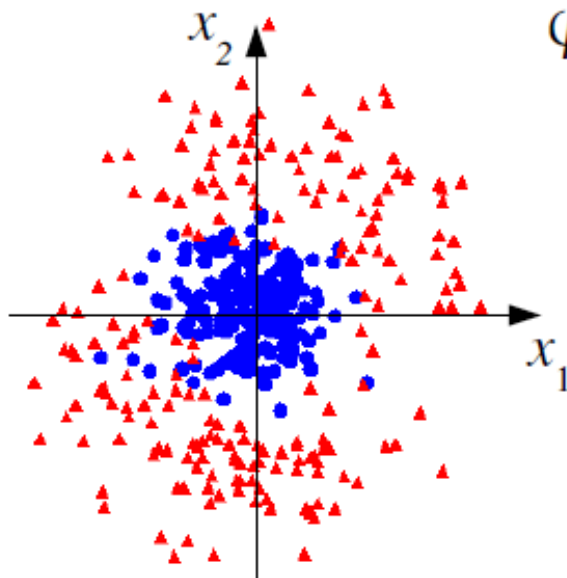
Nonlinear transformation of inputs

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:

$$\varphi_1 = \tan^{-1}(x_2/x_1)$$

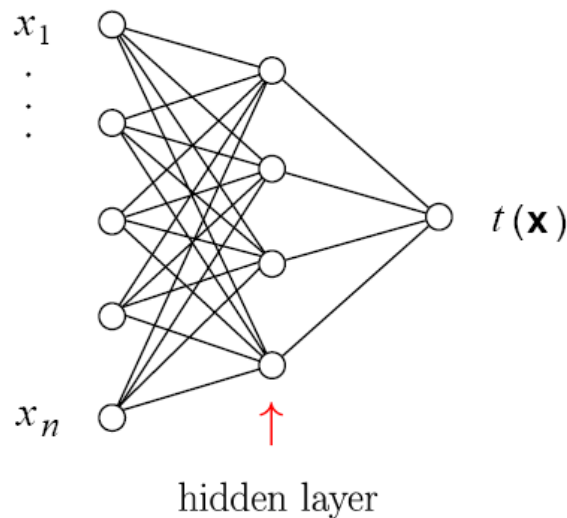
$$\varphi_2 = \sqrt{x_1^2 + x_2^2}$$

Here, guess fixed
basis functions
(no free parameters)



Neural networks in particle physics

For many years, the only "advanced" classifier used in particle physics.

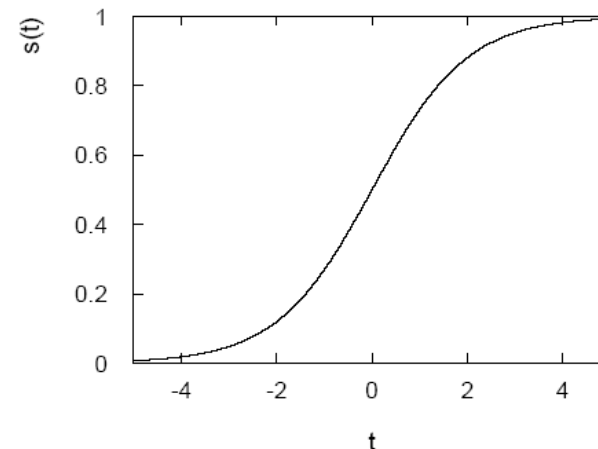


$$h_i(\vec{x}) = s \left(w_{i0} + \sum_{j=1}^n w_{ij} x_j \right) ,$$

$$t(\vec{x}) = s \left(a_0 + \sum_{i=1}^n a_i h_i(\vec{x}) \right) .$$

Usually use single hidden layer,
logistic sigmoid activation function:

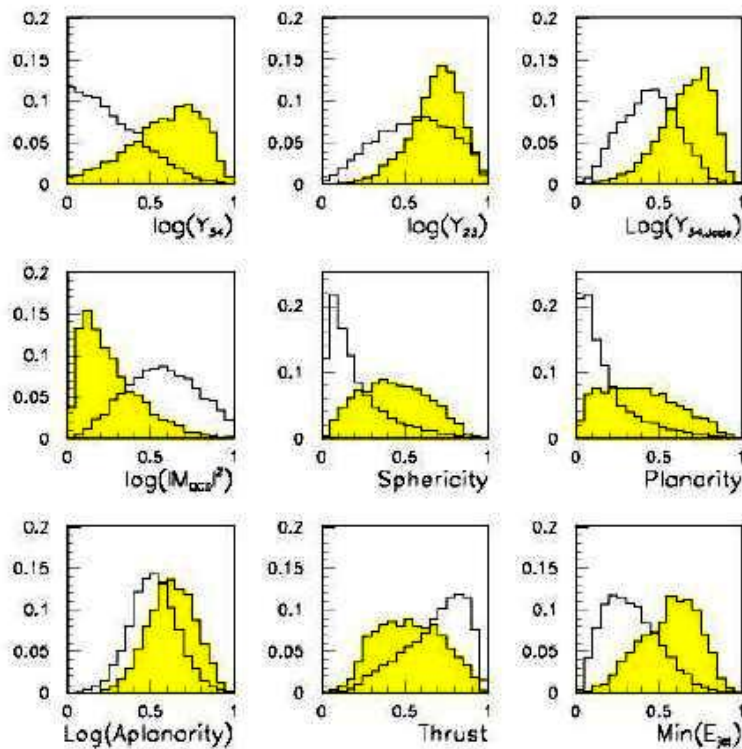
$$s(u) = (1 + e^{-u})^{-1}$$



Neural network example from LEP II

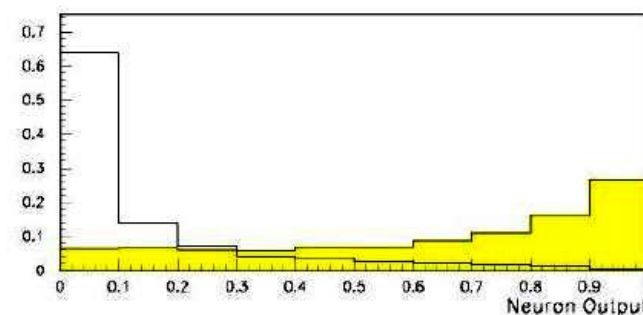
Signal: $e^+e^- \rightarrow W^+W^-$ (often 4 well separated hadron jets)

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



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

Neural network output:

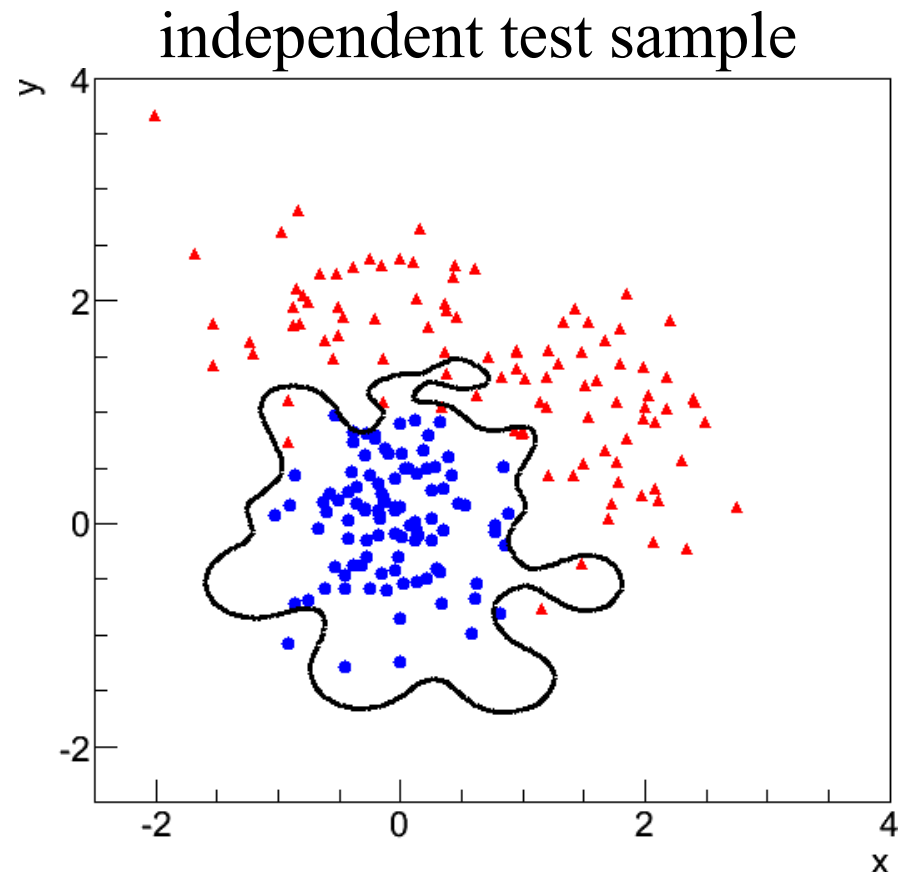
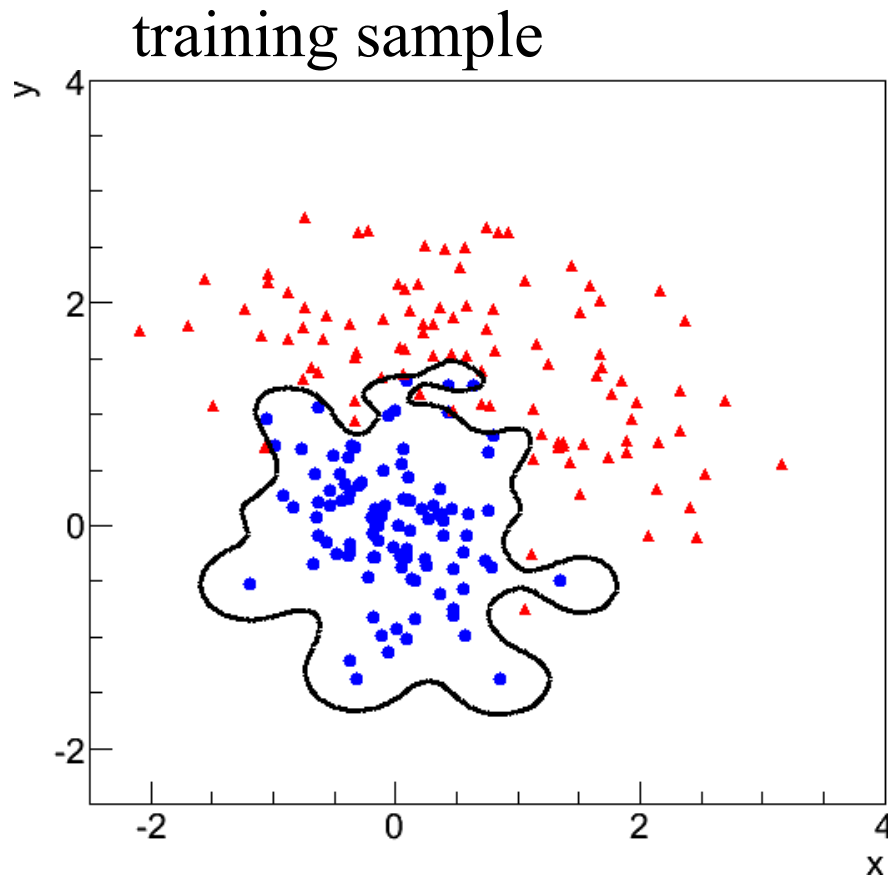


(Garrido, Juste and Martinez, ALEPH 96-144)

Overtraining

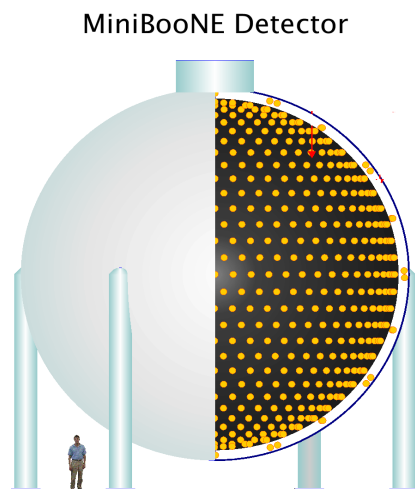
If decision boundary is too flexible it will conform too closely to the training points → **overtraining**.

Monitor by applying classifier to independent test sample.



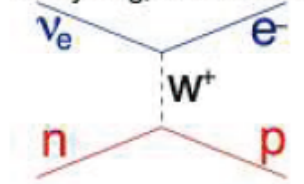
Particle i.d. in MiniBooNE

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

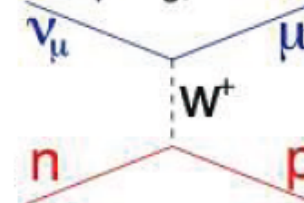


Search for ν_μ to ν_e oscillations required particle i.d. using information from the PMTs.

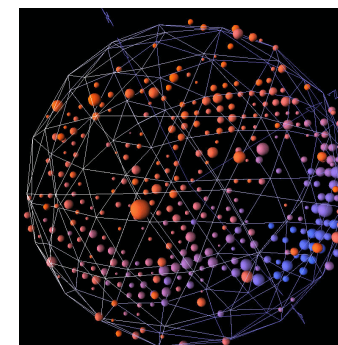
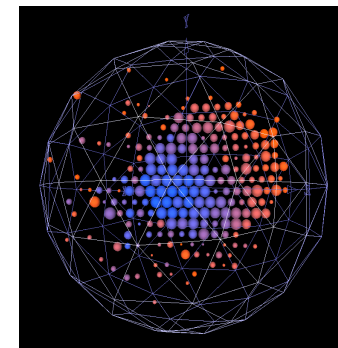
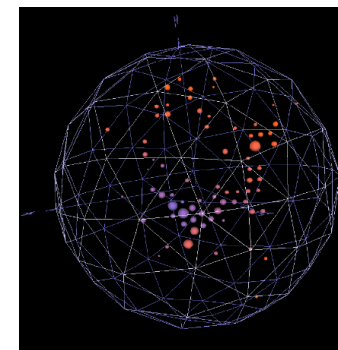
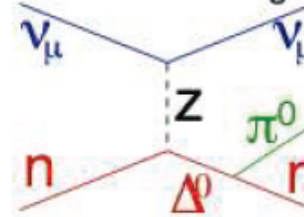
Electron candidate
fuzzy ring, short track



Muon candidate
sharp ring, filled in



Pion candidate
two "e-like" rings



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:

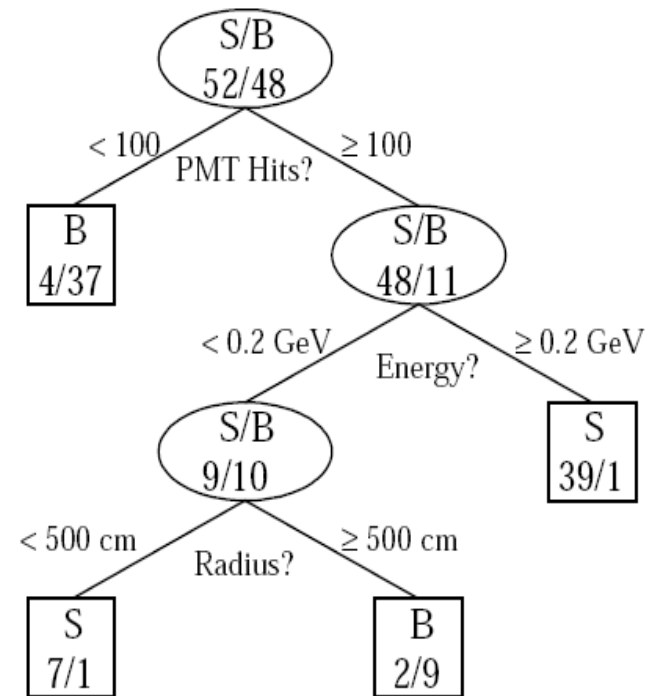
$$P = \frac{\sum_{\text{signal}} w_i}{\sum_{\text{signal}} w_i + \sum_{\text{background}} w_i}$$

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 Δ ; 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 as 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 the discriminant function

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

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.

AdaBoost

A successful boosting algorithm is AdaBoost (Freund & Schapire, 1997).

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

Train the classifier $f_1(\mathbf{x})$ (e.g. a decision tree) using the weights $w^{(1)}$ so as to minimize the classification error rate,

$$\varepsilon_1 = \sum_{i=1}^N w_i^{(1)} I(y_i f_1(\mathbf{x}_i) \leq 0),$$

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

Updating the event weights (AdaBoost)

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

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

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

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

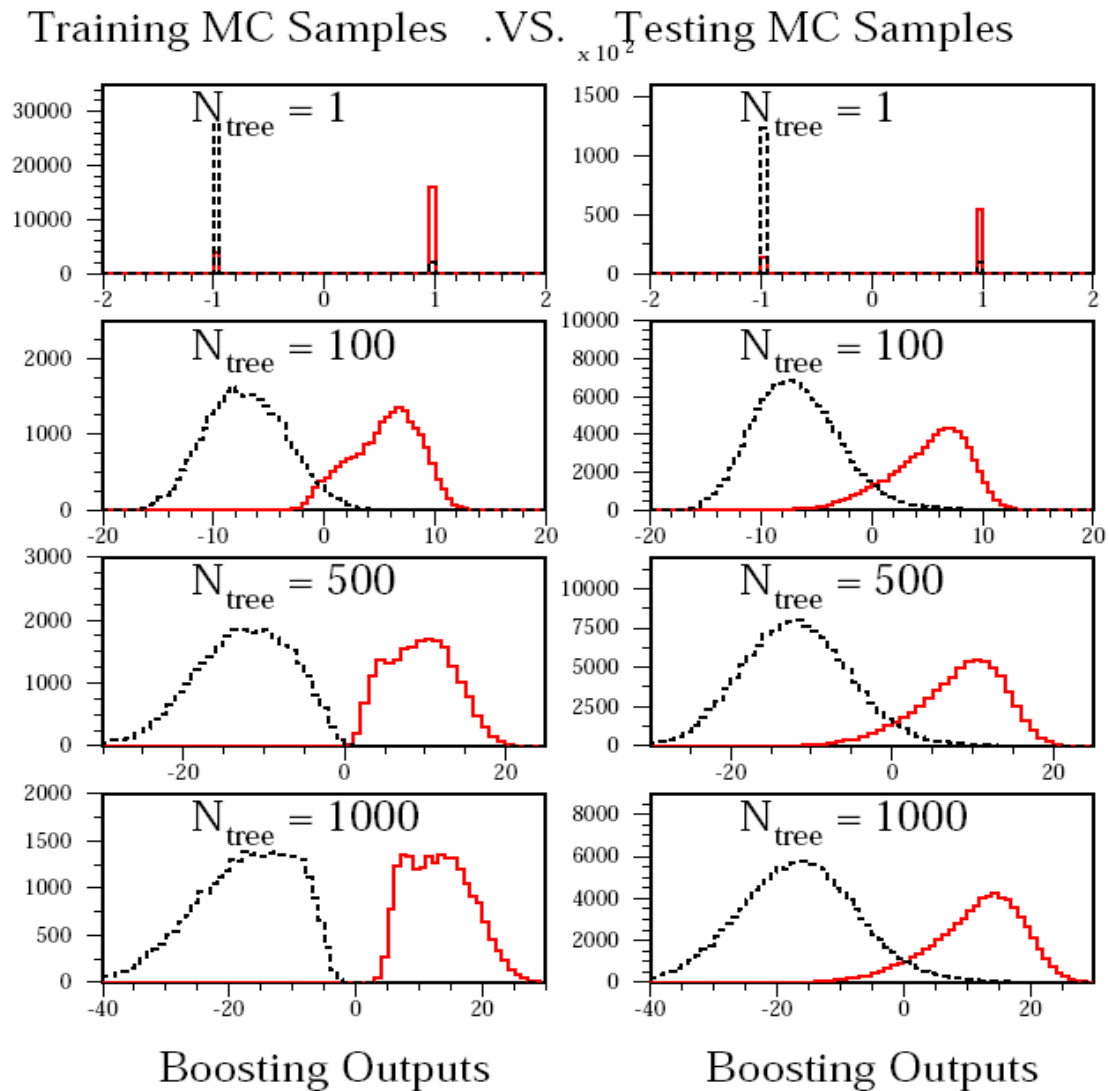
$i = \text{event index}$ $k = \text{training sample index}$ Normalize so that $\sum_i w_i^{(k+1)} = 1$

Iterate K times, final classifier is $y(\mathbf{x}) = \sum_{k=1}^K \alpha_k f_k(\mathbf{x}, T_k)$

Monitoring overtraining

From MiniBooNE
example:

Performance stable
after a few hundred
trees.



Boosted decision tree summary

Advantage of boosted decision tree is it can handle a large number of inputs. Those that provide little/no separation are rarely used as tree splitters are effectively ignored.

Easy to deal with inputs of mixed types (real, integer, categorical...).

If a tree has only a few leaves it is easy to visualize (but rarely use only a single tree).

There are a number of boosting algorithms, which differ primarily in the rule for updating the weights (ϵ -Boost, LogitBoost,...)

Other ways of combining weaker classifiers: Bagging (Bootstrap-Aggregating), generates the ensemble of classifiers by random sampling with replacement from the full training sample.

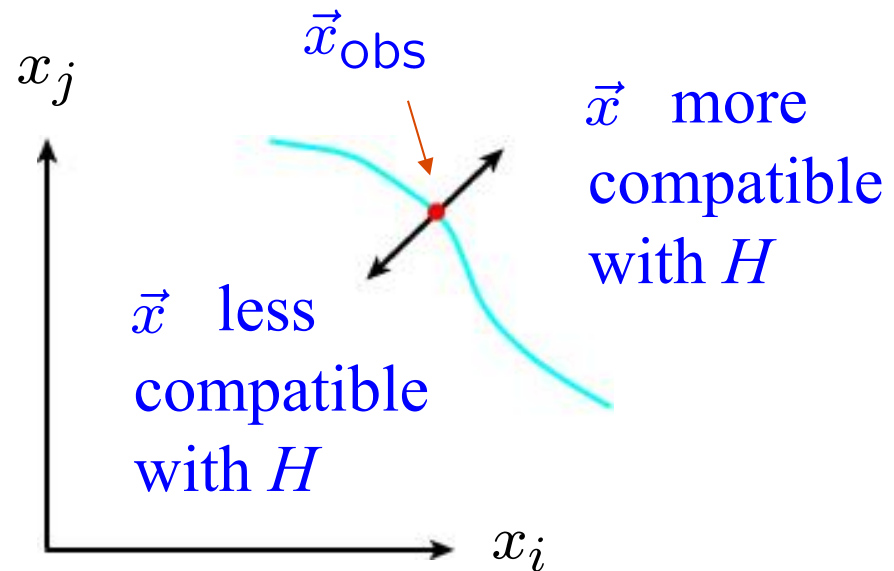
Testing significance / goodness-of-fit

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}_{obs}

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

Decide what part of the data space represents less compatibility with H than does the point \vec{x}_{obs} .
(Not unique!)



p-values

Express level of agreement between data and H with p -value:

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!

In frequentist statistics we don't talk about $P(H)$ (unless H represents a repeatable observation). In Bayesian statistics we do; use Bayes' theorem to obtain

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

where $\pi(H)$ is the prior probability for H .

For now stick with the frequentist approach;
result is p -value, regrettably easy to misinterpret as $P(H)$.

p-value example: testing whether a coin is ‘fair’

Probability to observe n heads in N coin tosses is binomial:

$$P(n; p, N) = \frac{N!}{n!(N - n)!} p^n (1 - p)^{N - n}$$

Hypothesis H : the coin is fair ($p = 0.5$).

Suppose we toss the coin $N = 20$ times and get $n = 17$ heads.

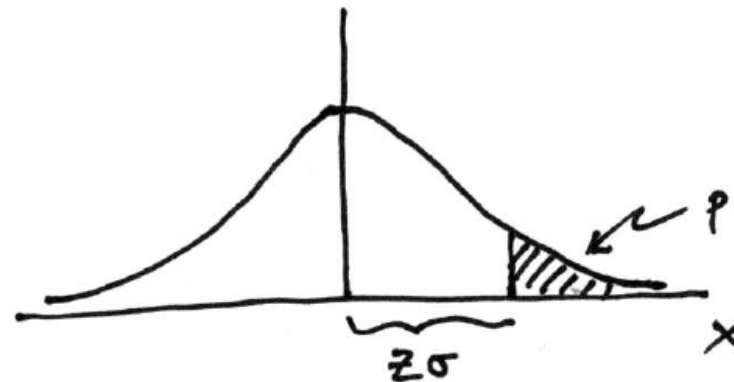
Region of data space with equal or lesser compatibility with H relative to $n = 17$ is: $n = 17, 18, 19, 20, 0, 1, 2, 3$. Adding up the probabilities for these values gives:

$$P(n = 0, 1, 2, 3, 17, 18, 19, \text{ or } 20) = 0.0026 .$$

i.e. $p = 0.0026$ is the probability of obtaining such a bizarre result (or more so) ‘by chance’, under the assumption of H .

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.



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

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

The significance of an observed signal

Suppose we observe n events; these can consist of:

n_b events from known processes (background)

n_s events from a new process (signal)

If n_s, n_b are Poisson r.v.s with means s, b , then $n = n_s + n_b$ is also Poisson, mean = $s + b$:

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

Suppose $b = 0.5$, and we observe $n_{\text{obs}} = 5$. Should we claim evidence for a new discovery?

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

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

Searching for presence of signal events

For each event we measure two variables, $\mathbf{x} = (x_1, x_2)$.

Suppose that for background events (hypothesis H_0),

$$f(\mathbf{x}|H_0) = \frac{1}{\xi_1} e^{-x_1/\xi_1} \frac{1}{\xi_2} e^{-x_2/\xi_2}$$

and for a certain signal model (hypothesis H_1) they follow

$$f(\mathbf{x}|H_1) = C \frac{1}{\sqrt{2\pi}\sigma_1} e^{-(x_1-\mu_1)^2/2\sigma_1^2} \frac{1}{\sqrt{2\pi}\sigma_2} e^{-(x_2-\mu_2)^2/2\sigma_2^2}$$

where $x_1, x_2 \geq 0$ and C is a normalization constant.

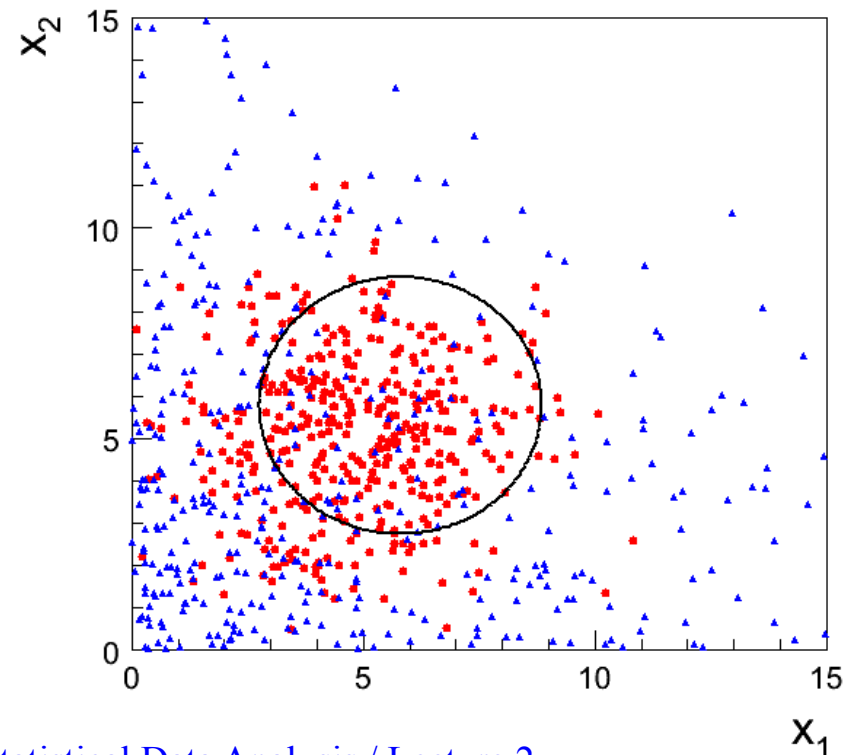
Likelihood ratio as test statistic

In a real-world problem we usually wouldn't have the pdfs $f(\mathbf{x}|H_0)$ and $f(\mathbf{x}|H_1)$, so we wouldn't be able to evaluate the likelihood ratio

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

for a given observed \mathbf{x} , hence the need for multivariate methods to approximate this with some other function.

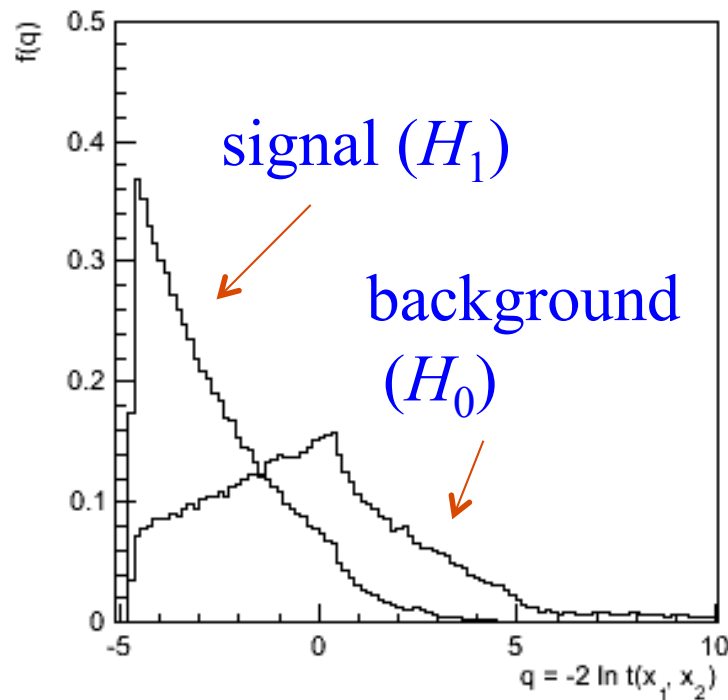
But in this example we can find contours of constant likelihood ratio such as:



Event selection using the LR

Using Monte Carlo, we can find the distribution of the likelihood ratio or equivalently of

$$q = \left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 - \frac{2x_1}{\xi_1} - \frac{2x_2}{\xi_2} = -2 \ln t(\mathbf{x}) + C$$



From the Neyman-Pearson lemma we know that by cutting on this variable we would select a signal sample with the highest signal efficiency (test power) for a given background efficiency.

Search for the signal process

But what if the signal process is not known to exist and we want to search for it. The relevant hypotheses are therefore

H_0 : all events are of the background type

H_1 : the events are a mixture of signal and background

Rejecting H_0 with $Z > 5$ constitutes “discovering” new physics.

Suppose that for a given integrated luminosity, the expected number of signal events is s , and for background b .

The observed number of events n will follow a Poisson distribution:

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

Likelihoods for full experiment

We observe n events, and thus measure n instances of $\mathbf{x} = (x_1, x_2)$.

The likelihood function for the entire experiment assuming the background-only hypothesis (H_0) is

$$L_b = \frac{b^n}{n!} e^{-b} \prod_{i=1}^n f(\mathbf{x}_i | b)$$

and for the “signal plus background” hypothesis (H_1) it is

$$L_{s+b} = \frac{(s+b)^n}{n!} e^{-(s+b)} \prod_{i=1}^n (\pi_s f(\mathbf{x}_i | s) + \pi_b f(\mathbf{x}_i | b))$$

where π_s and π_b are the (prior) probabilities for an event to be signal or background, respectively.

Likelihood ratio for full experiment

We can define a test statistic Q monotonic in the likelihood ratio as

$$Q = -2 \ln \frac{L_{s+b}}{L_b} = -s + \sum_{i=1}^n \ln \left(1 + \frac{s}{b} \frac{f(\mathbf{x}_i|s)}{f(\mathbf{x}_i|b)} \right)$$

To compute p -values for the b and $s+b$ hypotheses given an observed value of Q we need the distributions $f(Q|b)$ and $f(Q|s+b)$.

Note that the term $-s$ in front is a constant and can be dropped.

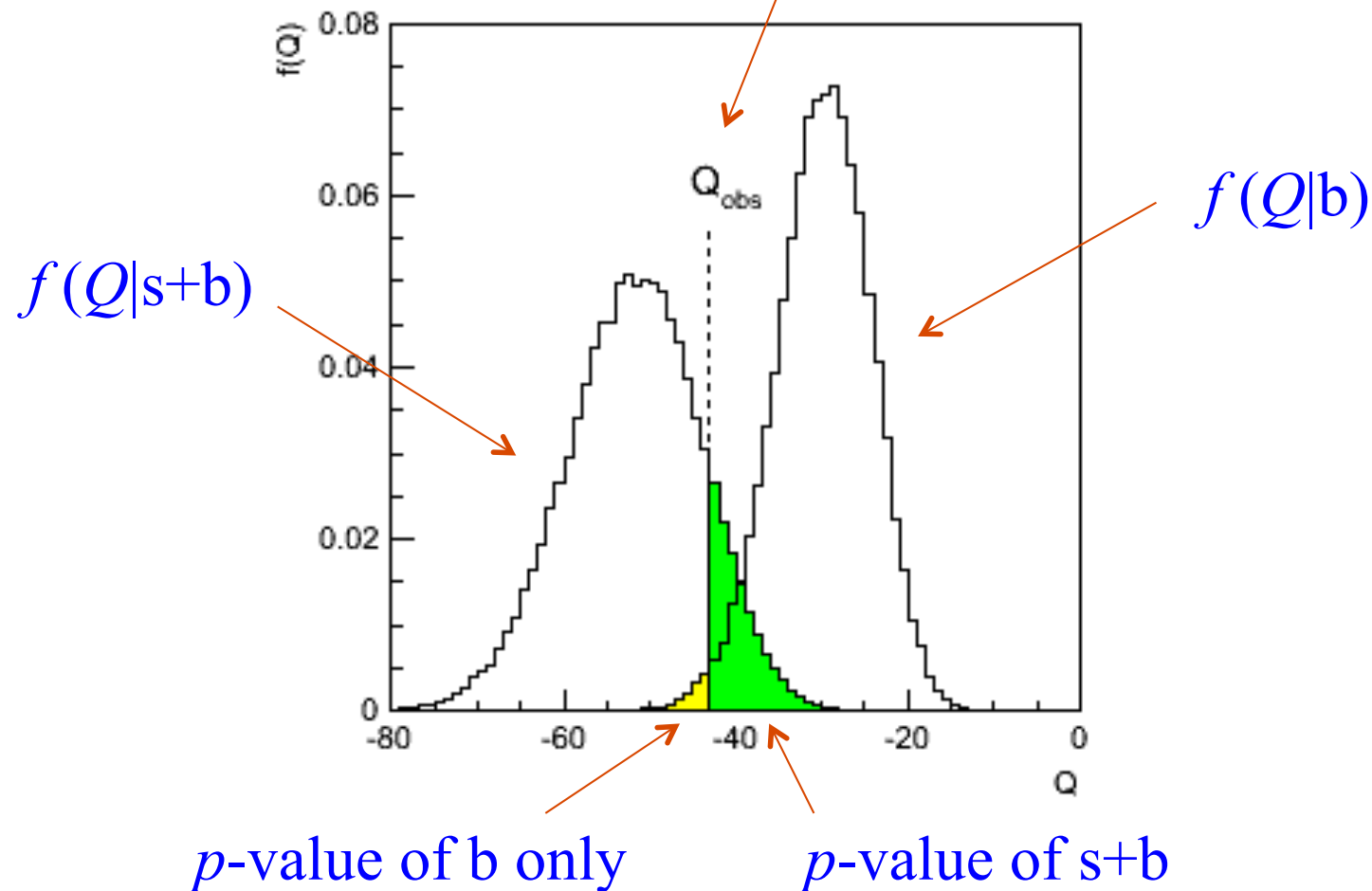
The rest is a sum of contributions for each event, and each term in the sum has the same distribution.

Can exploit this to relate distribution of Q to that of single event terms using (Fast) Fourier Transforms (Hu and Nielsen, physics/9906010).

Distribution of Q

Take e.g. $b = 100, s = 20$.

Suppose in real experiment
 Q is observed here.



Systematic uncertainties

Up to now we assumed all parameters were known exactly.

In practice they have some (systematic) uncertainty.

Suppose e.g. uncertainty in expected number of background events b is characterized by a (Bayesian) pdf $\pi(b)$.

Maybe take a Gaussian, i.e.,

$$\pi(b) = \frac{1}{\sqrt{2\pi}\sigma_b} e^{-(b-b_0)^2/2\sigma_b^2}$$

where b_0 is the nominal (measured) value and σ_b is the estimated uncertainty.

In fact for many systematics a Gaussian pdf is hard to defend – more on this later.

Distribution of Q with systematics

To get the desired p -values we need the pdf $f(Q)$, but this depends on b , which we don't know exactly.

But we can obtain the **Bayesian model average**:

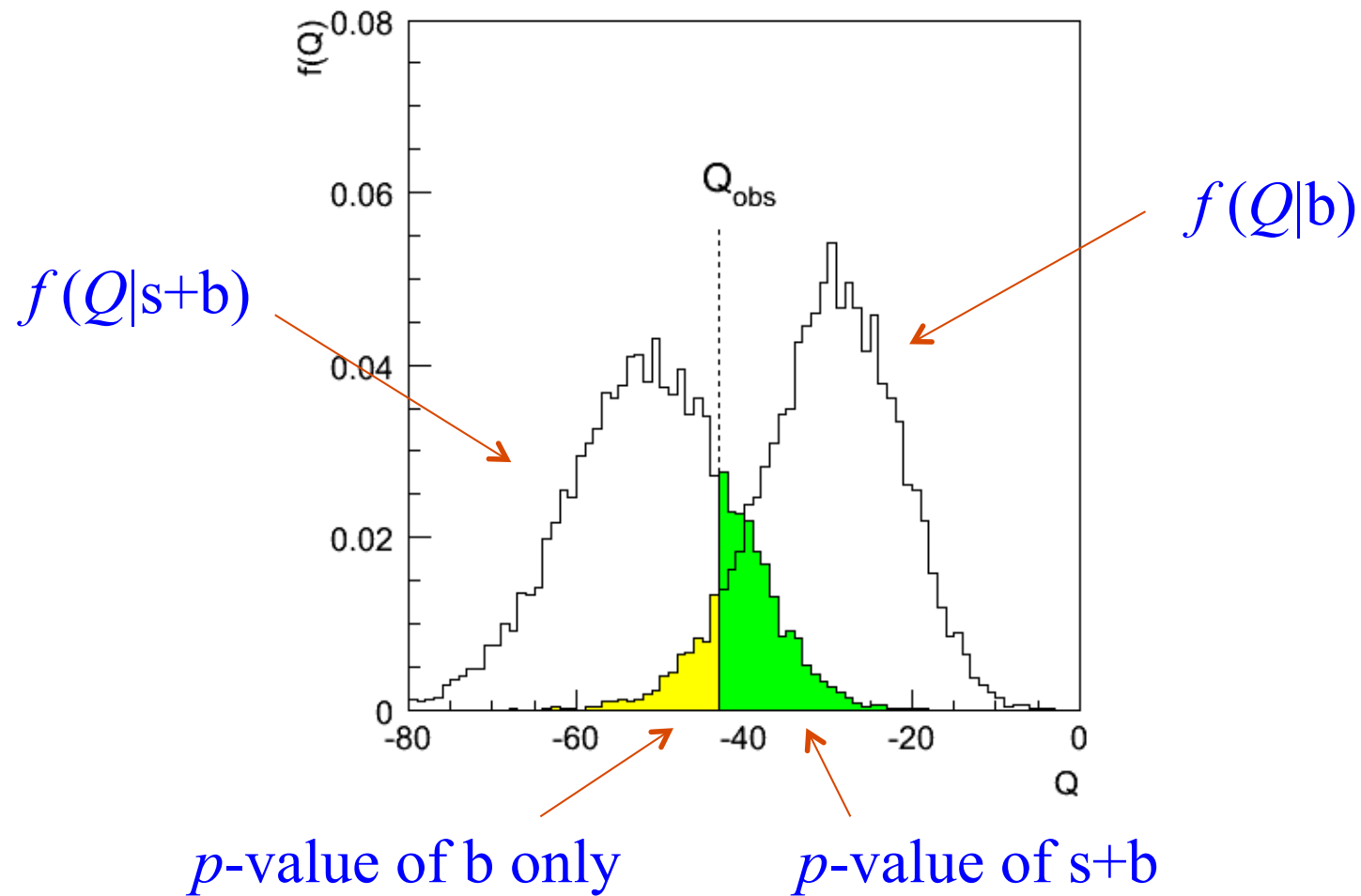
$$f(Q) = \int f(Q|b)\pi(b) db$$

With Monte Carlo, sample b from $\pi(b)$, then use this to generate Q from $f(Q|b)$, i.e., a new value of b is used to generate the data for every simulation of the experiment.

This broadens the distributions of Q and thus increases the p -value (decreases significance Z) for a given Q_{obs} .

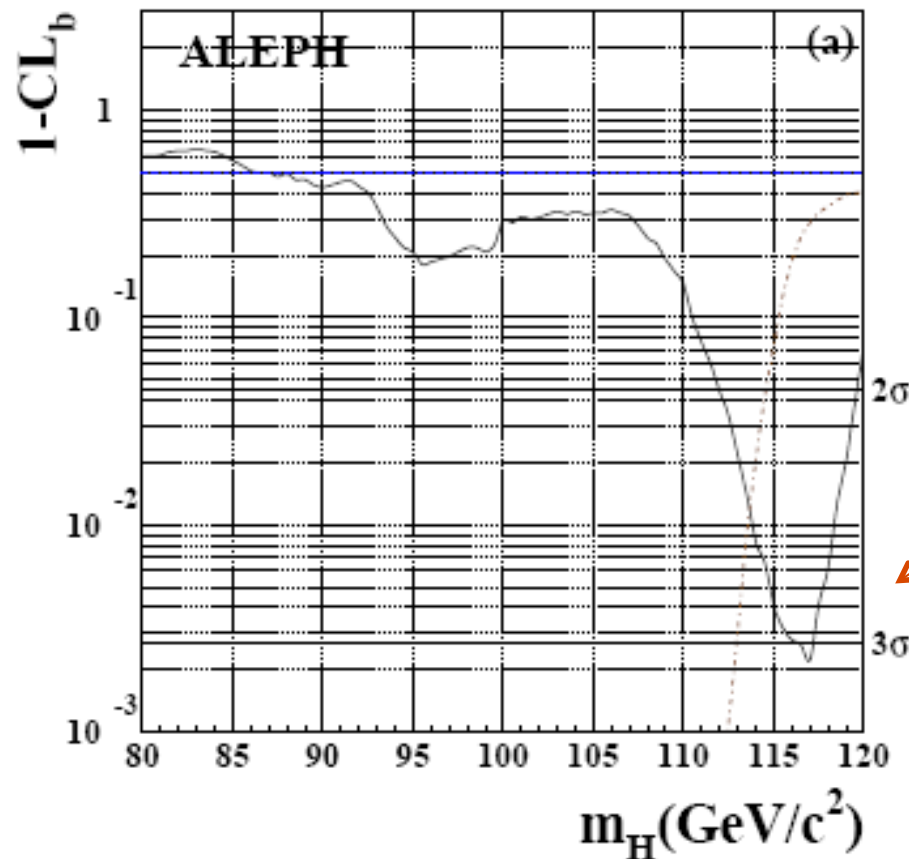
Distribution of Q with systematics (2)

For $s = 20$, $b_0 = 100$, $\sigma_b = 10$ this gives



Example: ALEPH Higgs search

p -value ($1 - \text{CL}_b$) of background only hypothesis versus tested Higgs mass measured by ALEPH Experiment

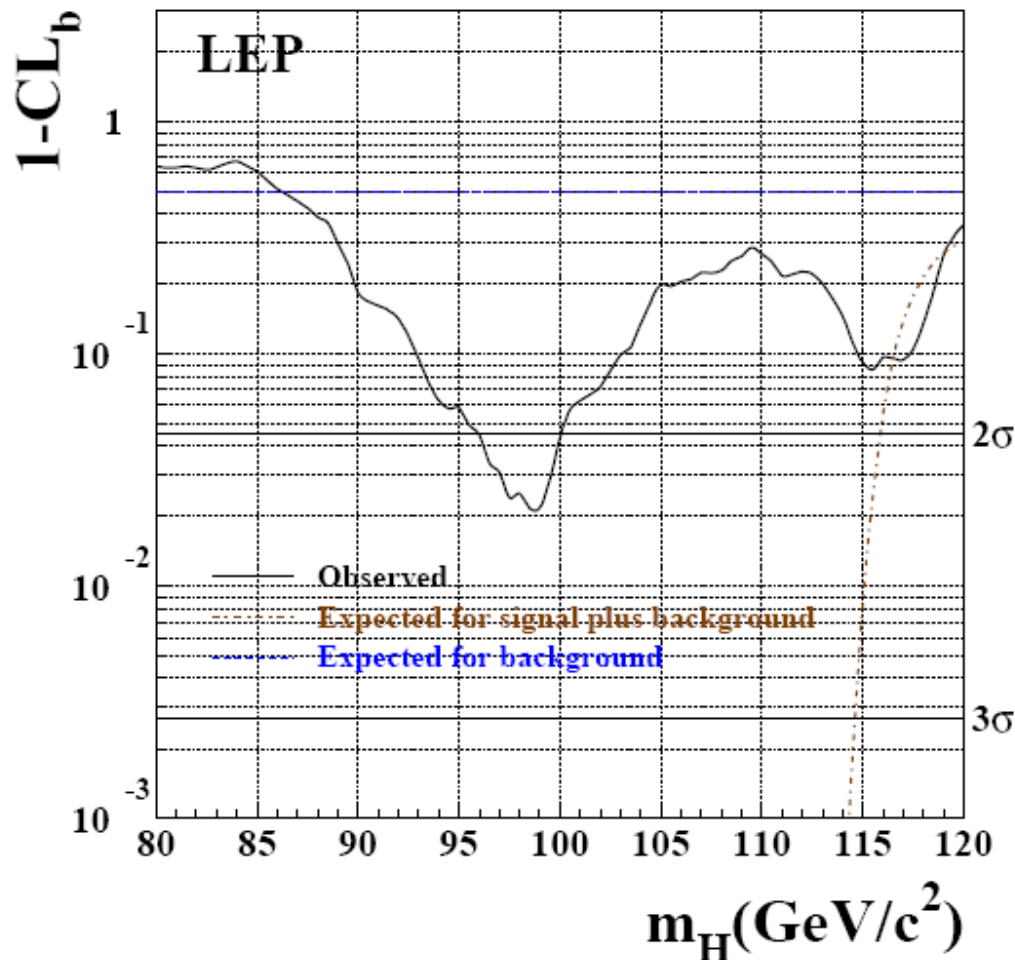


Possible signal?

Phys.Lett.B565:61-75,2003.
hep-ex/0306033

Example: LEP Higgs search

Not seen by the other LEP experiments. Combined analysis gives p -value of background-only hypothesis of 0.09 for $m_H = 115$ GeV.



Phys.Lett.B565:61-75,2003.
hep-ex/0306033

Using the likelihood ratio $L(s)/L(\hat{s})$

Instead of the likelihood ratio L_{s+b}/L_b , suppose we use as a test statistic

$$\lambda(s) = \frac{L(s)}{L(\hat{s})}$$



maximizes $L(s)$

Intuitively this is a good measure of the level of agreement between the data and the hypothesized value of s .

low λ : poor agreement

high λ : good agreement

$$0 \leq \lambda \leq 1$$

$L(s)/L(\hat{s})$ for counting experiment

Consider an experiment where we only count n events with $n \sim \text{Poisson}(s + b)$. Then $\hat{s} = n - b$.

To establish discovery of signal we test the hypothesis $s = 0$ using

$$\ln \lambda(0) = n \ln(b) - b - n \ln n + n$$

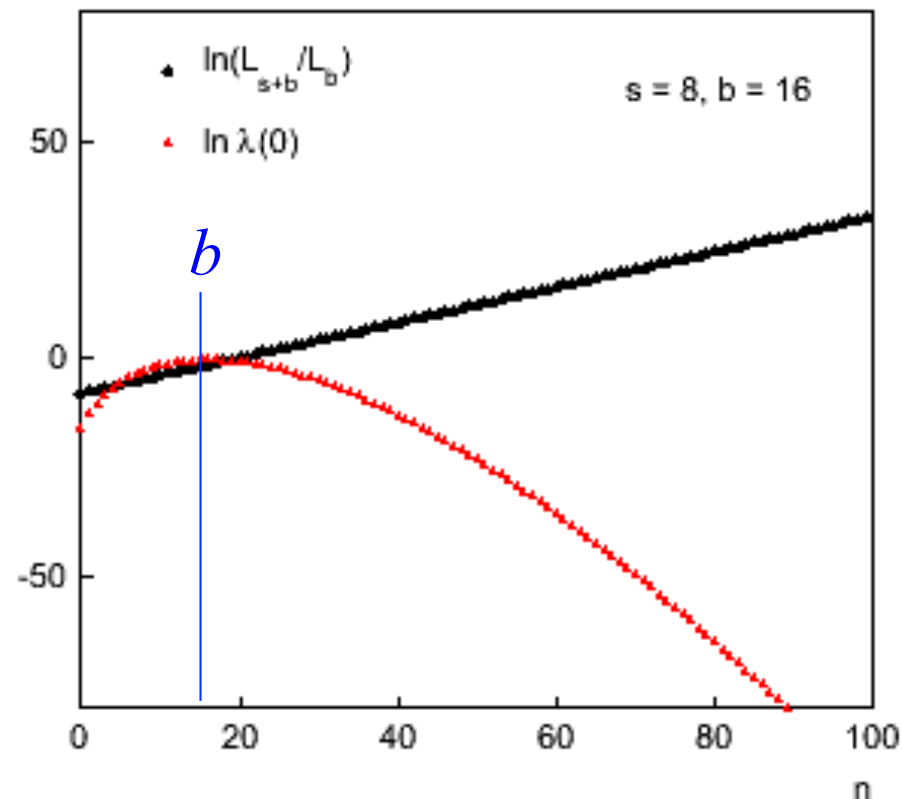
whereas previously we had used

$$\ln \frac{L_{s+b}}{L_b} = n \ln \left(1 + \frac{s}{b} \right) - s$$

which is monotonic in n and thus equivalent to using n as the test statistic.

$L(s)/L(\hat{s})$ for counting experiment (2)

But if we only consider the possibility of signal being present when $n > b$, then in this range $\lambda(0)$ is also monotonic in n , so both likelihood ratios lead to the same test.



$L(s)/L(\hat{s})$ for general experiment

If we do not simply count events but also measure for each some set of numbers, then the two likelihood ratios do not necessarily give equivalent tests, but in practice will be very close.

$\lambda(s)$ has the important advantage that for a sufficiently large event sample, its distribution approaches a well defined form (Wilks' Theorem).

In practice the approach to the asymptotic form is rapid and one obtains a good approximation even for relatively small data samples (but need to check with MC).

This remains true even when we have adjustable **nuisance parameters** in the problem, i.e., parameters that are needed for a correct description of the data but are otherwise not of interest (key to dealing with systematic uncertainties).

The profile likelihood ratio

If the model contains nuisance parameters θ , can base significance test on profile likelihood ratio:

$$\lambda(\mu) = \frac{L(\mu, \hat{\theta})}{L(\hat{\mu}, \hat{\theta})}$$

maximizes L for specified μ

maximize L

In large-sample limit, distribution of $-2\ln\lambda(\mu)$ related to chi-square pdf; see, e.g., Cowan, Cranmer, Gross, Vitells, EPJC 71 (2011) 1-19; arXiv:1007.1727,

Discovery significance for $n \sim \text{Poisson}(s + b)$

Consider again the case where we observe n events ,
model as following Poisson distribution with mean $s + b$
(assume b is known).

- 1) For an observed n , what is the significance Z_0 with which we would reject the $s = 0$ hypothesis?
- 2) What is the expected (or more precisely, median) Z_0 if the true value of the signal rate is s ?

Gaussian approximation for Poisson significance

For large $s + b$, $n \rightarrow x \sim \text{Gaussian}(\mu, \sigma)$, $\mu = s + b$, $\sigma = \sqrt{s + b}$.

For observed value x_{obs} , p -value of $s = 0$ is $\text{Prob}(x > x_{\text{obs}} \mid s = 0)$,:

$$p_0 = 1 - \Phi\left(\frac{x_{\text{obs}} - b}{\sqrt{b}}\right)$$

Significance for rejecting $s = 0$ is therefore

$$Z_0 = \Phi^{-1}(1 - p_0) = \frac{x_{\text{obs}} - b}{\sqrt{b}}$$

Expected (median) significance assuming signal rate s is

$$\text{median}[Z_0 \mid s + b] = \frac{s}{\sqrt{b}}$$

Better approximation for Poisson significance

Likelihood function for parameter s is

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

or equivalently the log-likelihood is

$$\ln L(s) = n \ln(s+b) - (s+b) - \ln n!$$

Find the maximum by setting $\frac{\partial \ln L}{\partial s} = 0$

gives the estimator for s : $\hat{s} = n - b$

Approximate Poisson significance (continued)

The likelihood ratio statistic for testing $s = 0$ is

$$q_0 = -2 \ln \frac{L(0)}{L(\hat{s})} = 2 \left(n \ln \frac{n}{b} + b - n \right) \quad \text{for } n > b, \quad 0 \text{ otherwise}$$

For sufficiently large $s + b$, (use Wilks' theorem),

$$Z_0 \approx \sqrt{q_0} = \sqrt{2 \left(n \ln \frac{n}{b} + b - n \right)} \quad \text{for } n > b, \quad 0 \text{ otherwise}$$

To find $\text{median}[Z_0|s+b]$, let $n \rightarrow s + b$,

$$\text{median}[Z_0|s + b] \approx \sqrt{2 \left((s + b) \ln(1 + s/b) - s \right)}$$

This reduces to s/\sqrt{b} for $s \ll b$.

Lecture 2 Summary

Neyman-Pearson lemma: likelihood ratio is optimal test statistic.
But usually not directly usable; try e.g.,

- Fisher discriminant

- Neural networks

- Boosted Decision Trees

- Support Vector Machines, ...

Significance tests

p -value of H is probability to see data with equal or worse compatibility with H (not same as $P(H)$).

“Discovery” = p -value of background-only hypothesis v. low

Systematic uncertainties

Quantified via nuisance parameters; methods include
Bayesian averaging, profile likelihood

Extra slides

Resources on multivariate methods

Books:

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

T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning*, Springer, 2001

R. Duda, P. Hart, D. Stork, *Pattern Classification*, 2nd ed., Wiley, 2001

A. Webb, *Statistical Pattern Recognition*, 2nd ed., Wiley, 2002

Materials from some recent meetings:

PHYSTAT conference series (2002, 2003, 2005, 2007,...) see
www.phystat.org

Caltech workshop on multivariate analysis, 11 February, 2008
indico.cern.ch/conferenceDisplay.py?confId=27385

SLAC Lectures on Machine Learning by Ilya Narsky (2006)
www-group.slac.stanford.edu/sluo/Lectures/Stat2006_Lectures.html

Software for multivariate analysis

TMVA, Höcker, Stelzer, Tegenfeldt, Voss, Voss, [physics/0703039](#)

From **tmva.sourceforge.net**, also distributed with ROOT

Variety of classifiers

Good manual

StatPatternRecognition, I. Narsky, [physics/0507143](#)

Further info from [www.hep.caltech.edu/~narsky/spr.html](#)

Also wide variety of methods, many complementary to **TMVA**

Currently appears project no longer to be supported