# Computing and Statistical Data Analysis
# Stat 4: MC, Intro to Statistical Tests

London Postgraduate Lectures on Particle Physics;
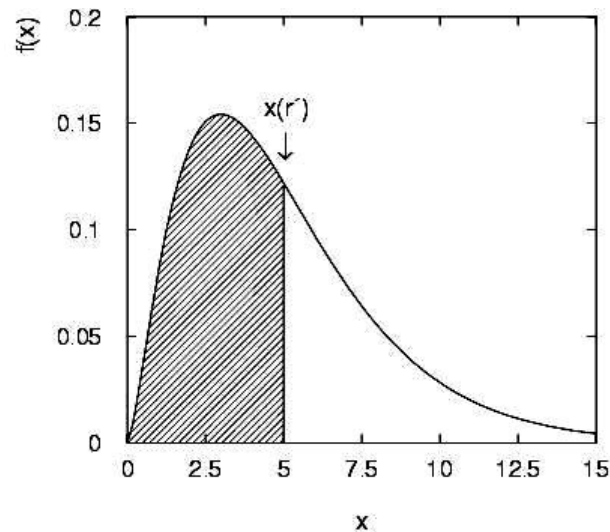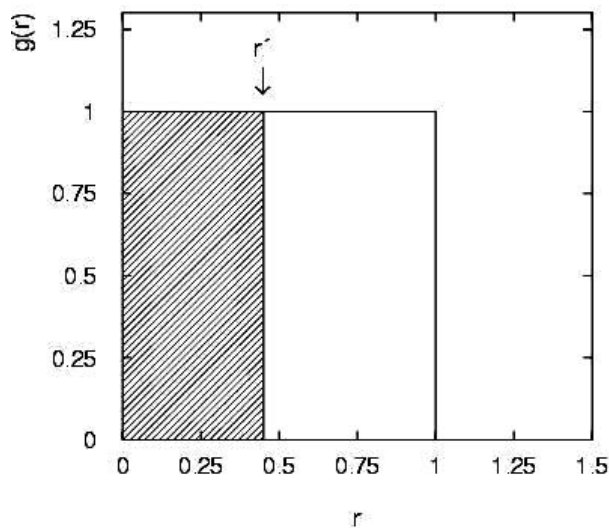
University of London MSci course PH4515

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

Course web page:

`www.pp.rhul.ac.uk/~cowan/stat_course.html`

# The transformation method

Given $r_1, r_2, ..., r_n$ uniform in $[0, 1]$, find $x_1, x_2, ..., x_n$
that follow $f(x)$ by finding a suitable transformation $x(r)$.



Require:  $P(r \leq r') = P(x \leq x(r'))$

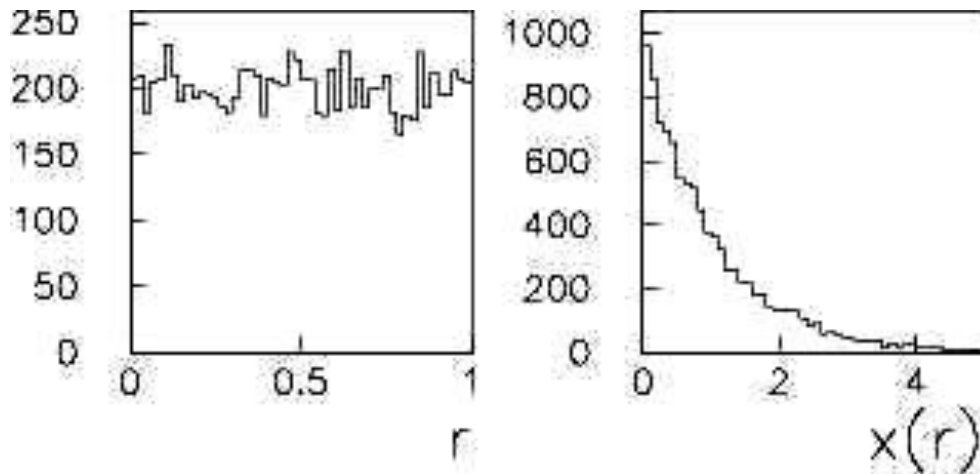i.e.  $\int_{-\infty}^{r'} g(r)\, dr = r' = \int_{-\infty}^{x(r')} f(x')\, dx' = F(x(r'))$

That is,  set $F(x) = r$ and solve for $x(r)$.

# Example of the transformation method

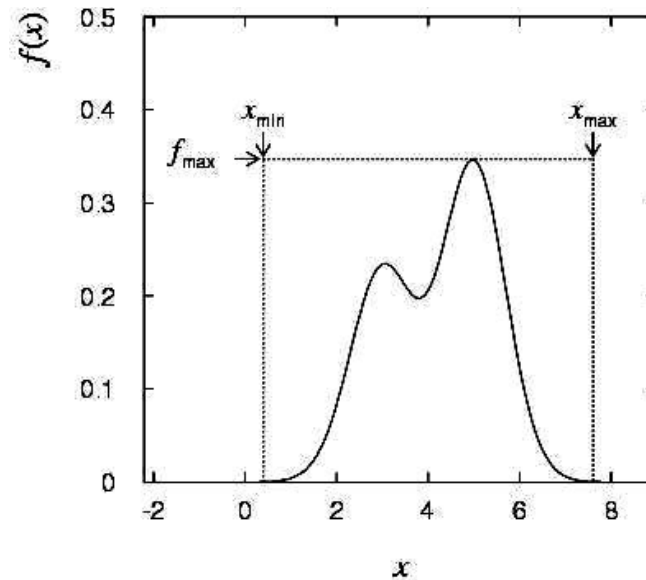Exponential pdf: $\quad f(x; \xi) = \dfrac{1}{\xi} e^{-x/\xi} \quad (x \geq 0)$

Set $\quad \displaystyle\int_0^x \frac{1}{\xi} e^{-x'/\xi}\, dx' = r \quad$ and solve for $x\,(r)$.

$\longrightarrow \quad x(r) = -\xi \ln(1 - r) \quad (\, x(r) = -\xi \ln r \,$ works too.$)$
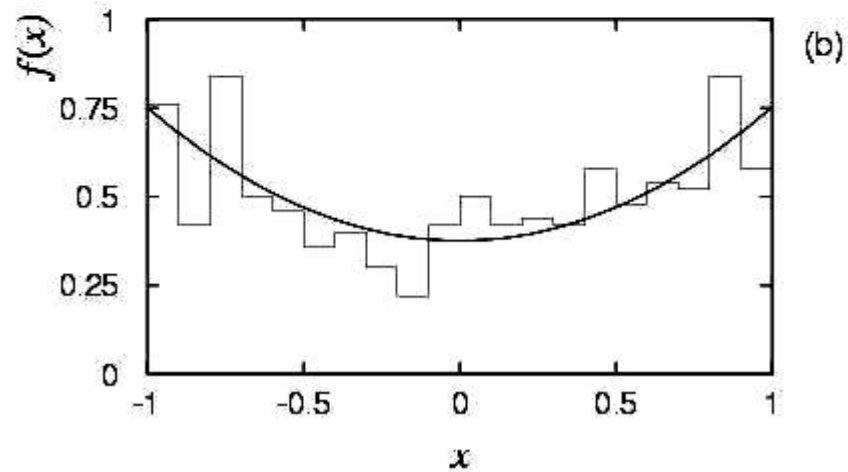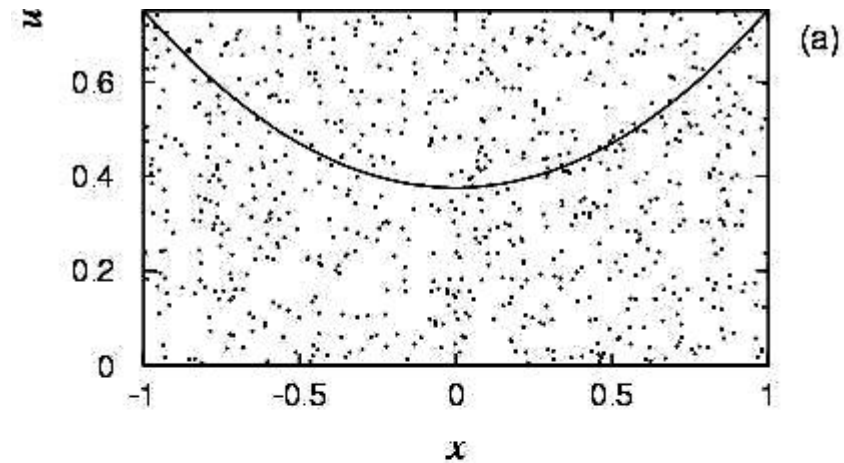
# The acceptance-rejection method

Enclose the pdf in a box:



(1) Generate a random number $x$, uniform in $[x_{min}, x_{max}]$, i.e.
$x = x_{min} + r_1(x_{max} - x_{min})$ , $r_1$ is uniform in [0,1].

(2) Generate a 2nd independent random number $u$ uniformly distributed between 0 and $f_{max}$, i.e. $u = r_2 f_{max}$ .

(3) If $u < f(x)$, then accept $x$. If not, reject $x$ and repeat.

# Example with acceptance-rejection method

$$f(x) = \frac{3}{8}(1 + x^2)$$

$$(-1 \leq x \leq 1)$$

If dot below curve, use $x$ value in histogram.

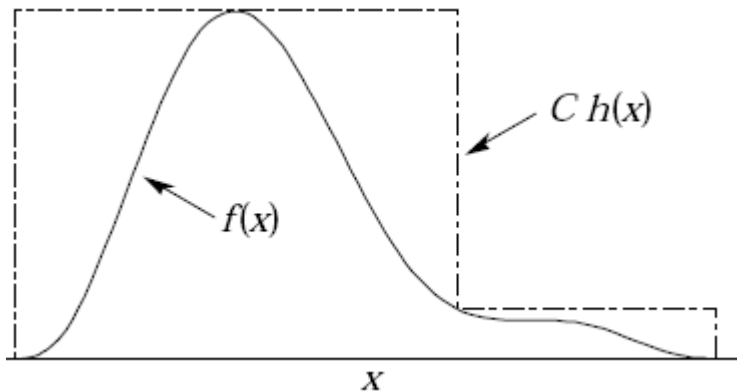Computing and Statistical Data Analysis / Stat 4

# Improving efficiency of the acceptance-rejection method

The fraction of accepted points is equal to the fraction of the box's area under the curve.

For very peaked distributions, this may be very low and thus the algorithm may be slow.

Improve by enclosing the pdf $f(x)$ in a curve $C\,h(x)$ that conforms to $f(x)$ more closely, where $h(x)$ is a pdf from which we can generate random values and $C$ is a constant.
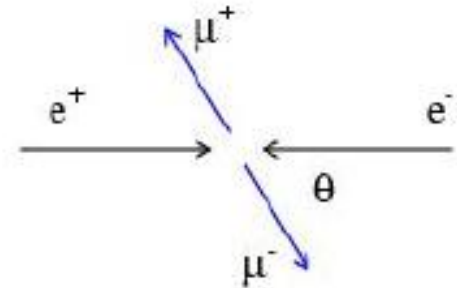


Generate points uniformly over $C\,h(x)$.

If point is below $f(x)$, accept $x$.

# Monte Carlo event generators

Simple example:  $e^+e^- \to \mu^+\mu^-$

Generate $\cos\theta$ and $\phi$:

$$f(\cos\theta; A_{\mathsf{FB}}) \propto (1 + \frac{8}{3}A_{\mathsf{FB}}\cos\theta + \cos^2\theta) \,,$$

$$g(\phi) = \frac{1}{2\pi} \quad (0 \le \phi \le 2\pi)$$

Less simple:  'event generators' for a variety of reactions:

$e^+e^- \to \mu^+\mu^-$, hadrons, ...

$pp \to$ hadrons, D-Y, SUSY,...

e.g. PYTHIA, HERWIG, ISAJET...

Output = 'events', i.e., for each event we get a list of generated particles and their momentum vectors, types, etc.

# A simulated event



PYTHIA Monte Carlo
pp → gluino-gluino

# Monte Carlo detector simulation

Takes as input the particle list and momenta from generator.

Simulates detector response:

multiple Coulomb scattering (generate scattering angle),
particle decays (generate lifetime),
ionization energy loss (generate $\Delta$),
electromagnetic, hadronic showers,
production of signals, electronics response, ...

Output = simulated raw data → input to reconstruction software:
track finding, fitting, etc.

Predict what you should see at 'detector level' given a certain hypothesis for 'generator level'. Compare with the real data.

Estimate 'efficiencies' = #events found / # events generated.

Programming package: `GEANT`

# 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.,

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

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:

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

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

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

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

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

# 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

Muon Detectors    Tile Calorimeter    Liquid Argon Calorimeter

Toroid Magnets    Solenoid Magnet    SCT Tracker    Pixel Detector    TRT Tracker

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

# A simulated SUSY event

high $p_T$ muons

high $p_T$ jets of hadrons

ATLAS   Atlantis   Event: susyevent

p

p

missing transverse energy

# Background events



This event from Standard Model ttbar production also has high $p_\text{T}$ jets and muons, and some missing transverse energy.

$\rightarrow$ 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 $\vec{x} = (x_1, \ldots, x_n)$

$x_1$ = 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
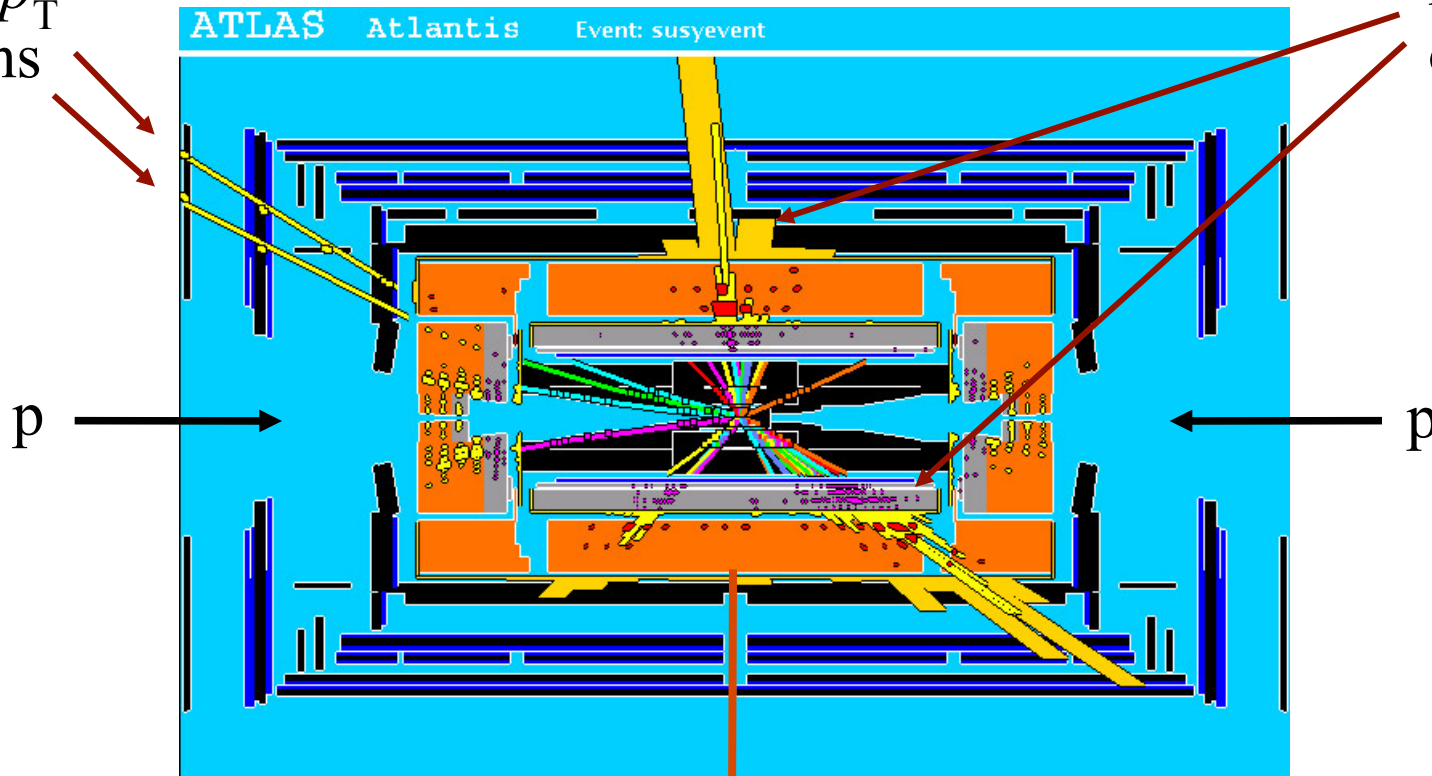
$$\text{pp} \rightarrow t\bar{t} , \qquad \text{pp} \rightarrow \tilde{g}\tilde{g} , \ldots$$

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':

$$x_i \quad < c_i$$

$$x_j \quad < c_j$$

# 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 decision boundary can be defined by an equation of the form

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

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

We can work out the pdfs $\quad g(t|H_0), \; g(t|H_1), \; \ldots$

Decision boundary is now a single 'cut' on $t$, which divides the space into the critical (rejection) region and acceptance region.

This defines a test. If the data fall in the critical region, we reject $H_0$.

# Signal/background efficiency

Probability to reject background hypothesis for background event (background efficiency):

$$\varepsilon_{\mathrm{b}} = \int_{t_{\mathrm{cut}}}^{\infty} g(t|\mathrm{b})\, dt = \alpha$$

Probability to accept a signal event as signal (signal efficiency):

$$\varepsilon_{\mathrm{s}} = \int_{t_{\mathrm{cut}}}^{\infty} g(t|\mathrm{s})\, dt = 1 - \beta$$

# Purity of event selection

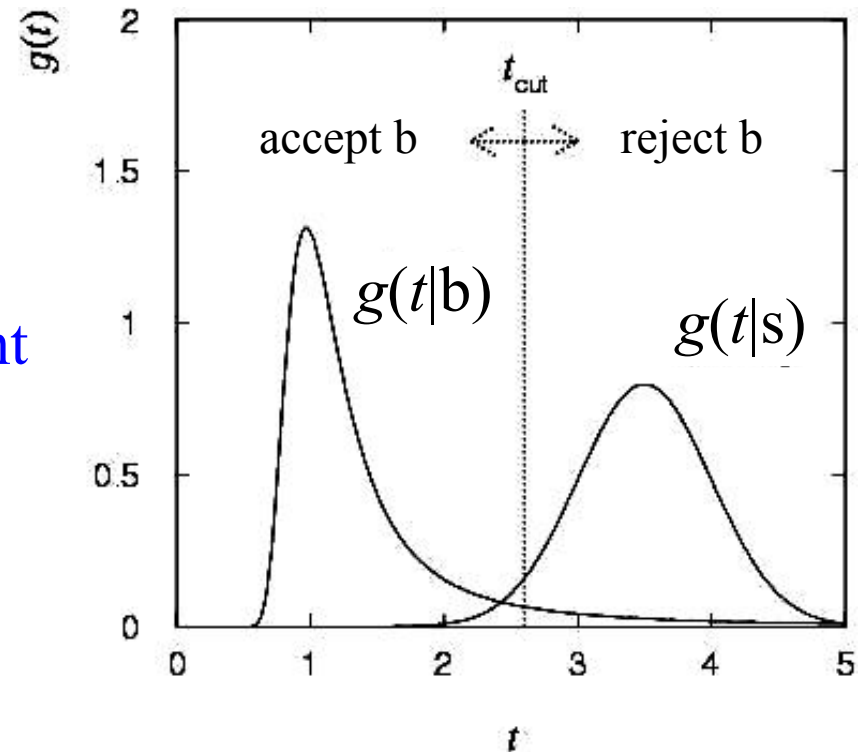Suppose only one background type b; overall fractions of signal and background events are $\pi_s$ and $\pi_b$ (prior probabilities).

Suppose we select signal events with $t > t_{cut}$. What is the 'purity' of our selected sample?

Here purity means the probability to be signal given that the event was accepted. Using Bayes' theorem we find:

$$P(s|t > t_{cut}) = \frac{P(t > t_{cut}|s)\pi_s}{P(t > t_{cut}|s)\pi_s + P(t > t_{cut}|b)\pi_b}$$

$$= \frac{\varepsilon_s \pi_s}{\varepsilon_s \pi_s + \varepsilon_b \pi_b}$$

So the purity depends on the prior probabilities as well as on the signal and background efficiencies.

# Constructing a test statistic

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

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

inside the region, and $\leq c$ outside, where $c$ is a constant which determines the power.

Equivalently, optimal scalar test statistic is $\boxed{t(\mathbf{x}) = \frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)}}$

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

# Why Neyman-Pearson doesn't always help

The problem is that we usually don't have explicit formulae for the pdfs $P(x|H_0)$, $P(x|H_1)$.

Instead we may have Monte Carlo models for signal and background processes, so we can produce simulated data, and enter each event into an $n$-dimensional histogram.

Use e.g. $M$ bins for each of the $n$ dimensions, total of $M^n$ cells.

But $n$ is potentially large, $\rightarrow$ prohibitively large number of cells to populate with Monte Carlo data.

Compromise: make Ansatz for form of test statistic $t(\vec{x})$ with fewer parameters; determine them (e.g. using MC) to give best discrimination between signal and background.

# 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