

#### **Monte Carlo methods**



## The Monte Carlo method

A numerical technique for calculating probabilities & related things using sequences of random numbers.

## The usual steps:

- generate sequence  $r_1, r_2, ..., r_m$  uniform in ]0,1[.
- use them to produce another sequence  $x_1, x_2, ..., x_n$  distributed according to pdf f(x) of interrest
- use obtained x values to estimate some property of

f(x), e.g. fraction of x values within  $[a,b] = \int_a^b f(x) dx$ 

 $\Rightarrow$  MC calculation sort of integration (at least formally) Trivial for 1D:  $\int_a^b f(x) dx$  obtainable by other methods, but MC more powerful for multi-dimensional problems.

# MC x values = "simulated data" → used for testing statistical procedures

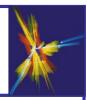
MC methods a wide field, actually own field in itself – here focus on the usage of MC for data analysis e.g. determining the statistical (& systematic) uncertainties. In such cases, MCs are used to generate different data distributions. So, let's try to answer the question "how can I generate the type of distribution I need?".

more complete & deeper discussion found in e.g.

• MATR323 "Basics of Monte Carlo simulation"-course by Prof. Flyura Djurabekova in spring 2026 (period III)



#### Random number generator



goal: to get uniformly distributed values in ]0,1[ interval. ⇒ "random number generator"

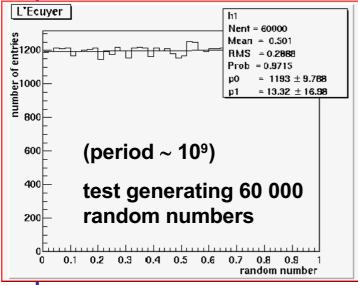
= computer algorithm to generate  $r_1, r_2, ..., r_m$ .

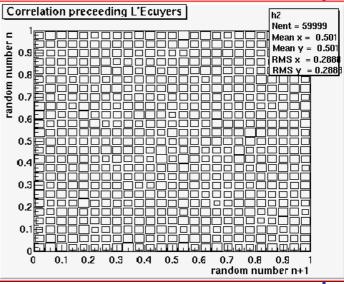
e.g. multiplicative linear congruential generator (MLCG)

$$n_{i+1} = (an_i) \mod m$$
, where

 $n_i$  = integer, a = **multiplier**, m = **modulus** &  $n_0$  = **seed**. (NB! mod = modulus (remainder), e.g. 27 mod 5 = 2)  $n_i$  follow periodic sequence in [1, m-1]  $\Rightarrow$   $r_i = n_i / m$  distributed in ]0,1[.

Choose a & m so that  $r_i$ 's pass various tests of randomness:  $r_i$ 's uniform in ]0,1[, succeeding  $r_i$ 's uncorrelated & "period" for  $r_i$ 's long (maximum = m –1) e.g. L'Ecuyer, Comm. ACM 31(1988)742: a = 40692, m = 2147483399

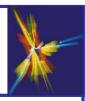




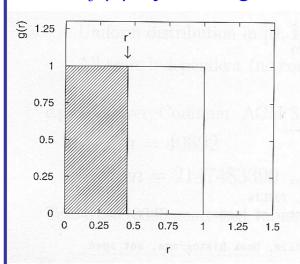
Much better algorithms exist e.g. Mersenne twister, with period  $\approx 10^{6000}$ . Many good algorithms implemented in freely available program libraries NB!  $r_i$ 's like above in reality **pseudorandom numbers** See e.g. F.James, *Comput. Phys. Commun.* 60 (1990) 111

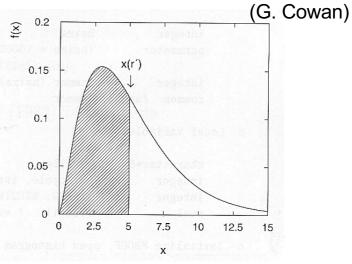


#### Inverse transform method



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





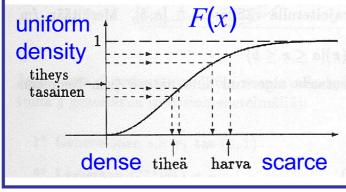
$$P(r \le r') = P(x \le x(r'))$$
 i.e.  

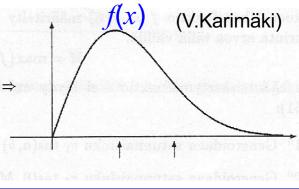
$$\int_{-\infty}^{r'} g(r)dr = r' = \int_{-\infty}^{x(r')} f(x')dx' = F(x(r'))$$

A method that always works when inverse function of cumulative distribution function F(x) can be calculated or put in a table. Then **inverse transform method**:

- sample r from a uniform distribution ]0,1[
- calculate  $x = F^{-1}(r)$

then generated random numbers x that obey pdf f(x). From following graph it is easy to see that it works





#### Inverse transform method



## **Discrete distributions**

Inverse transform method for discrete distributions.  $p_i$  = probability for integer i. First one has to put in table cumulative distribution function  $F_j = \sum_{i=0}^{j} p_i$ , j = 0, ..., N. If infinite number of possible outcomes i, N to be set so large that  $F_N \approx 1$ . Generation algorithm for discrete pdf:

- (i) sample *r* from a uniform distribution ]0,1[.
- (ii) find k so that  $F_{k-1} < r < F_k$ .
- (iii) accept integer k-1

Resulting distribution proportional to probabilities  $p_k$ .

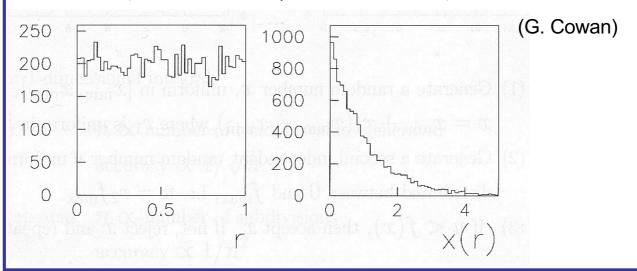
Example of inverse transform method:

exponential pdf: 
$$f(x;\xi) = \xi^{-1}e^{-x/\xi}$$
  $(x \ge 0)$ 

cumulative distribution function:

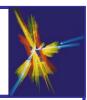
$$F(x) = \int_0^x \xi^{-1} e^{-x'/\xi} dx' = 1 - e^{-x/\xi}$$

assume  $r \in ]0,1[$ , now set r = F(x) & solve for  $x(r) \Rightarrow x(r) = -\xi \ln(1-r)$  (NB!  $x(r) = -\xi \ln r$  works also)



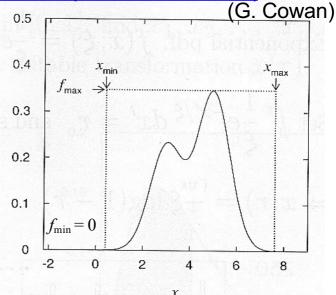


## UNIVERSITY OF HELSINKI Acceptance-rejection method



Acceptance-rejection method (von Neumann)

Often analytic solution  $\stackrel{\mathfrak{Z}}{\approx}$  impossible or very impractical  $\Rightarrow$  acceptance-rejection method (or **hit-or-miss**): enclose pdf in a box  $]f_{\min} = \min(f(x)), f_{\max} = \max(f(x))[$ 

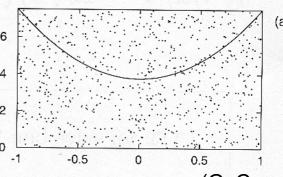


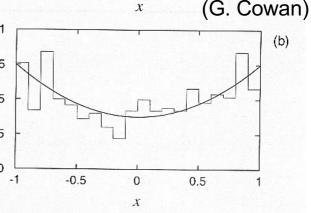
(i) generate a random number x, uniform in  $]x_{\min}$ ,  $x_{\max}[$ , i.e.  $x = x_{\min} + r_1(x_{\max} - x_{\min})$  where  $r_1$  uniform in ]0,1[ (ii) generate  $2^{\text{nd}}$  random number u uniformly distributed between  $f_{\min}$  &  $f_{\max}$ , i.e.  $u = f_{\min} + r_2(f_{\max} - f_{\min})$ . (iii) if u < f(x), then accept x. If not, reject x & repeat.

Example: generate a polar old angle distribution  $1 + \cos \theta^2$  old  $(-1 \le \cos \theta \le 1)$ .  $x = \cos \theta$ 

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

 $f_{\text{min}} = 0$ ,  $f_{\text{max}} = \frac{3}{4}$ ; x of pairs of random numbers lying below the curve accepted. 0.25 The distribution of the accepted x shown below.



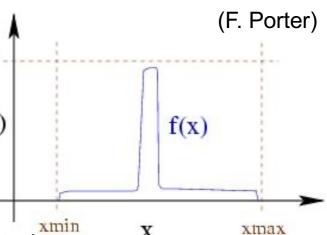




#### Importance sampling



Acceptance-rejection method simple to apply but efficiency of algorithm depends on f(x)area ratio of pdf to enclosing box. Inefficient for "peaky" distributions.



 $\int_{x_{min}}^{x_{max}} f(x) dx$ 

Fraction of trial points accepted:

$$\varepsilon = \frac{\text{Area under curve}}{\text{Area of box}} = \frac{\int_{x_{min}}^{x_{max}} f(x) dx}{(f_{max} - f_{min})(x_{max} - x_{min})}$$

## Importance sampling:

To improve acceptance-rejection method efficiency, use importance sampling method. First random numbers generated according to g(x) satisfying g(x) > f(x) all over whole x-interval. x's generated according to  $g(x)/\int g(x)$  (with e.g. inverse transform method), and x accepted if u < f(x), where u random number ]0, g(x)[, or weighted with a factor f(x)/g(x)

P(Y > 3) if Y Gaussian  $N(\mu = 0, \sigma = 1)$ : truth  $\approx$ e.g. to estimate the 0.00135.  $h(Y_i) = 1$  if  $Y_i > 3$ ;  $h(Y_i) = 0$  if  $Y_i \le 3$ 

integral of the tail of  $\circ$  Draw an iid sample  $Y_1, \ldots, Y_{100}$  from a N(0,1), then the

Gaussian distribution. estimator is (J. Cisewski)

Instead of using the original Gaussian,

use an equivalent

Gaussian with larger acceptance rate to reduce uncertainty

of the estimate and coverge much faster. integral  $\hat{I} = \frac{1}{100} \sum_{i=1}^{100} h(Y_i)$ 

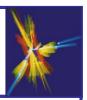
• Draw an iid sample  $Y_1, \ldots, Y_{100}$  from a N(4,1), then the estimator is

where f is the density of a N(0,1) and g is the density of N(4,1)

$N=10^5$	Expected Value	Variance
Truth	0.00135	0
Monte Carlo	0.00136	$1.3 \times 10^{-5}$
Importance Sampling	0.00135	$9.5 \times 10^{-8}$



#### MC in particle physics



MC methods often used to "mimic" data in e.g. particle physics, a two-step process: simulation of physics, "event generator" & simulation of response of experimental apparatus, "detector simulation".

MC event generator: simple example:  $e^+e^- \rightarrow \mu^+\mu^-$  generate  $\theta$  and  $\phi$ :  $f(\cos\theta; A_{FB}) \propto (1 + \frac{8}{3} A_{FB} \cos\theta + \cos^2\theta)$   $\mu^ g(\phi) = 1/2\pi$ 

In reality implemented into program packages that accounts for all (known) effects (hadronisation, initial and final state radiation, longlived particles etc...)

MC detector simulation (built on e.g. GEANT4): Input: particle list & momenta from event generator

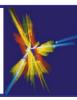
simulate particle interactions with detector material & detector response: multiple Coulomb scattering (generate particle scattering angle), ionization energy loss (generate energy loss dE/dx), electromagnetic & hadronic showers, produce detector signals, electronics response ...

Output: simulated "raw" data → input to reconstruction Usage:

Predict what should be seen at "detector level" given a hypothesis at "generator level". Compare with real data

- optimize measurement & experiment sensitivity.
- estimate efficiency & purity (expected signal & background)
- simulate measurement many times ("gedanken" experiment)





### **Standard Gaussian distributed random numbers**

cumulative distribution function 
$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-x^{-2}/2} dx'$$

Square evaluation method: Generate numbers in polar coordinates & afterwards transform them to cartesian

$$f(x)f(y)dxdy = re^{-r^2/2}dr\frac{d\phi}{2\pi} = g(r)h(\phi)drd\phi$$

 $r \& \phi$  generated according to  $g(r) \& h(\phi) \Rightarrow$ 

x & y independent Gaussian distributed random numbers.

Marsaglia polar method: very efficient algorithm for generating Gaussian random numbers

 $u_1 \& u_2$  uniform in ]0,1[

construct  $v_1 = 2u_1 - 1 \& v_2 = 2u_2 - 1$  (uniform in ]-1,1[),

if  $r^2 = v_1^2 + v_2^2 > 1$  start over again, otherwise

$$z_1 = v_1 \sqrt{-2 \ln r^2 / r^2}$$
 and  $z_2 = v_2 \sqrt{-2 \ln r^2 / r^2}$ 

 $z_1$  &  $z_2$  independent Gaussian numbers with  $\mu$  = 0 &  $\sigma$  = 1  $z_i$ ' =  $\mu$  +  $\sigma z_i$  Gaussian numbers with mean  $\mu$  & variance  $\sigma^2$ 

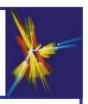
Poisson random numbers: efficient algorithm for small  $\nu$ : Set k = 1 & A = 1 at start, then iterate until successful choice

- (i) generate u (uniform in ]0,1[), replace A with uA
- (ii) if  $A < \exp(-\nu)$ , where  $\nu$  mean of Poisson distribution, accept  $n_k = k-1$  & stop.
- (iii) replace k by k+1 and repeat (i).

for large  $\nu$  (> ~10), faster & easier to generate Gaussian random numbers since Poisson distribution  $\approx$  Gaussian.



#### **Accuracy of MC methods**



 $\chi^2$  random numbers: for even n, generate n/2 uniform

random numbers 
$$u_i$$
; then  $y = -2 \ln \left( \prod_{i=1}^{n/2} u_i \right)$  follows  $\chi^2(n)$ .

for odd n, generate (n-1)/2 uniform numbers  $u_i$  & one

Gaussian z; then 
$$y = -2 \ln \left( \prod_{i=1}^{(n-1)/2} u_i \right) + z^2$$
 follows  $\chi^2(n)$ .

Binomial random numbers: principle same as for any discrete distribution. Use inverse transform method and tabulate cumulative distribution F(x) Most computer libraries include generators for most common distributions like Gaussian, Poisson,  $\chi^2$  etc...

## **Accuracy of Monte Carlo methods:**

MC = "integration". (G. Cowan) compare to trapezoidal rule, n = # of computing steps 0.15 for 1D integral: MC:  $n \propto$ number of accepted random 0.1 values, accuracy  $\propto 1/\sqrt{n}$ trapezoid:  $n \propto$  number of 0.05 subdivisions, accuracy  $\propto 1/n^2$ in 1D trapezoid wins!!! 2.5 7.5 12.5

MC: accuracy  $\propto 1/\sqrt{n}$ 

but in *d* dimensions:

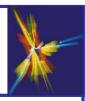
 $\leftarrow$  independent of d!

trapezoid: accuracy  $\propto 1/n^{2/d}$ 

MC wins for d > 4. Gaussian quadrature better than trapezoid but for high enough d, MC always wins!! (see e.g. F. James, *Rep. Prog. Phys.* 43 (1980) 1145).



#### **Hypothesis testing**



## A statistical test:

applied to physics mainly in two different ways

how to distinguish events of interest, "signal", from (a large number of) uninteresting events, "background"
how well a given hypothesis compatible (without any alternative) with observed data, "goodness-of-fit" test.

## **Hypothesis testing:**

Statistics related to questions where answers not numerical but logical – a "yes" or "no".

Example: "Is this particle an electron / a signal event?" Statistics (& science in general) can never in absolute terms say that something is true.

Only falsify ("empirical falsification" Karl Popper 1930).

Answer generally a probability.

This probability refers to ensemble of statements.

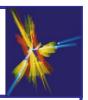
Result of a measurement  $\bar{x} = (x_1, ..., x_n)$  follows a pdf in n-dimensional space that depends on variables  $x_1, ..., x_n$  e.g. pdf  $f(\bar{x})$  specified by some **hypotheses**  $H_0, H_1...$  each having a probability density

$$f(\bar{x} \mid H_0), f(\bar{x} \mid H_1)$$
 etc..

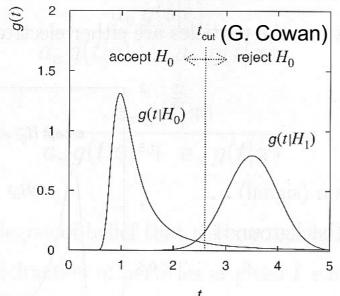
Equivalent to provide a statement & ask if true or not. null hypothesis:  $H_0$ , hypothesis whose validity tested. alternative hypothesis:  $H_1(, H_2 ...)$  to compare with  $H_0$ . simple hypothesis:  $f(\bar{x})$  completely specified. composite hypothesis:  $f(\bar{x}, \theta)$ , parameter(s)  $\theta$  unknown.



#### **Test statistic**



Many dimensions akward  $\Rightarrow$  use **test statistic**  $t(\overline{x})$  of lower dimension (e.g. 1D) to compactify data, keeping as much hypothesis discrimination power as possible. A test statistic t following pdf's  $g(t(\overline{x})|H_0), g(t(\overline{x})|H_1)$ 



Often formulate compatibility between data & various hypotheses in terms of an acceptence or a rejection of null hypothesis  $H_0$ 

**Critical region** i.e. where t not likely to occur if  $H_0$  true  $t \ge t_{\text{cut}}$  in figure (alternative define **acceptance region**)

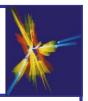
If observed value  $t_{\rm obs}$  in critical region, reject  $H_0$ , otherwise "accept" (or strictly speaking, not reject !!).

Critical region chosen such that probability for t to be there assuming hypothesis  $H_0$  some small value  $\alpha$  = significance level.  $\alpha = \int_{t_{cut}}^{\infty} g(t \mid H_0) dt$ 

 $H_0$  rejection even if true <u>error of 1<sup>st</sup> kind</u> (= probability  $\alpha$ )  $H_0$  acceptence if some other hypothesis  $t_{cut}$  true <u>error of 2<sup>nd</sup> kind</u> (probability  $\beta$ )  $\beta = \int g(t \mid H_1) dt$ 

Can define  $1-\beta$  as <u>power</u> to  $-\infty$  discriminate against alternative hypothesis.



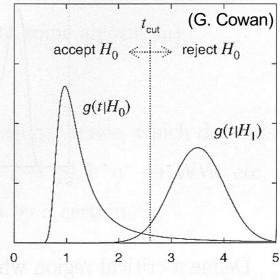


n energy loss measurements for a particle in a particle detector, t = truncated mean of the measurements & suppose all particles either electrons (e) or pions  $(\pi)$ .

 $H_0$  = electron (signal)  $\frac{2}{5}$   $t_1^2$  = pion (background)  $t_2^2$  select electrons with  $t_1^2$  <  $t_2^2$  selection efficiencies:

$$\varepsilon_{\rm e} = \int_{-\infty}^{t_{cut}} g(t \mid e) dt = 1 - \alpha \quad _{0.5}$$

$$\varepsilon_{\pi} = \int_{-\infty}^{t_{cut}} g(t \mid \pi) dt = \beta$$



Higher  $t_{\text{cut}}$ : higher e efficiency but more  $\pi$  background. Lower  $t_{\text{cut}}$ : lower e efficiency but better sample purity.

For observed value t, probability P, to have an e (or  $\pi$ )?

$$P(e \mid t) = \frac{a_{e}g(t \mid e)}{a_{e}g(t \mid e) + a_{\pi}g(t \mid \pi)} \quad P(\pi \mid t) = \frac{a_{\pi}g(t \mid \pi)}{a_{e}g(t \mid e) + a_{\pi}g(t \mid \pi)}$$

 $a_e(a_\pi)$  = fraction of electrons (pions) in sample;  $a_e+a_\pi$  = 1

Bayesian: degree of belief that particle is an e or a  $\pi$ . Frequentist: particle fraction at given t which is e or  $\pi$ . sample purity,  $p_{\rm e}$ , fraction of actual electrons in sample:

$$p_{e} = \int_{-\infty}^{t_{cut}} a_{e}g(t \mid e)dt / \int_{-\infty}^{t_{cut}} [a_{e}g(t \mid e) + a_{\pi}g(t \mid \pi)]dt$$

= electron probability averaged over interval ( $-\infty$ ,  $t_{\rm cut}$ ]. NB! purity depend on (un)known e &  $\pi$  fractions  $a_{\rm e}$  &  $a_{\pi}$ .



#### Discriminant analysis



## When is discriminant analysis needed?

aim of discriminant analysis: distinguish between a number of distinct hypotheses,  $H_k$ , k of  $m_k$ ,  $m_k$  various hypotheses of  $m_k$ ,  $m_k$ .

known a priori  $(H_1, H_2, ..., H_{n_h})$ .

problem: determine to which hypothesis new element (x)(determined a posteriori by simulation, measurement,...)

(determined a posteriori by simulation, measurement,...) belongs. Careful hypothesis  $H_1$  testing required, especially if only few elements available (low statistics).

 $\in H_{1}?, H_{2}?, ..., H_{n_{h}}?$ 

Discriminant analysis techniques applied in many fields e.g. astrophysics, particle physics, biophysics, imaging,...

A multidimensional test statistic  $\bar{t} = (t_1, ..., t_m)$  and hypotheses  $H_0$  ("signal") &  $H_1$  ("background"). Optimal choice of critical region, i.e. what selection to use?

Neyman-Pearson lemma: acceptance region giving

highest power (& also highest signal purity) for a given significance level  $\alpha$  (or selection efficiency  $\varepsilon = 1-\alpha$ ):

equivalently, optimal test statistic i.e. the **likelihood ratio** for simple shypotheses  $H_0 \& H_1$ . Requiring r > c gives maximum significance level (efficiency) for a given power (purity).

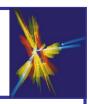
$$\frac{g(\bar{t}\mid H_0)}{g(\bar{t}\mid H_1)} > c,$$

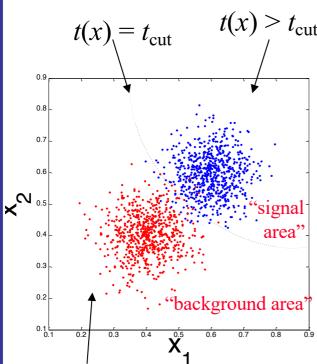
c = constant, fixed by significance level (or selection purity)

$$r = \frac{g(\bar{t} \mid H_0)}{g(\bar{t} \mid H_1)},$$



#### Constructing a test statistic





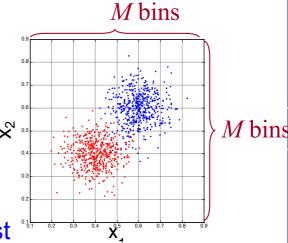
 $t(x) > t_{\text{cut}}$  A vector  $\bar{x} = (x_1, ..., x_n)$ measured for each event. construct 1D test statistic  $t(\bar{x})$  to distinguish between two hypotheses  $H_0 \& H_1$ . choice: the  $t(\bar{x}) = \frac{f(\bar{x} \mid H_0)}{f(\bar{x} \mid H_0)}$ most optimal likelihood ratio but need pdfs analytically  $f(\overline{x} \mid H_0) \& f(\overline{x} \mid H_1)$ (in most cases not possible).

In practice calculated, from MC simulations,  $t(x) < t_{\text{cut}}$  where pdf's approximated by multidimensional histograms filled with  $\bar{x}$  of each generate event.

method impractical if number of input variables too large since available MC statistics finite.

compromise: make Ansatz for functional form of  $t(\bar{x})$ with a fewer number of input variables; choose variables (e.g. based on MC) that give best discrimination between  $H_0 \& H_1$ .

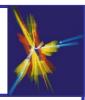
linear test statistic: e.g. Fisher discriminant - straight "line" cut in  $\bar{x}$ -space:  $t(\bar{x}) = \sum_{i=1}^{n} a_i x_i = \bar{a}^T \bar{x}$  more complicated cut function in  $\bar{x}$ -space.



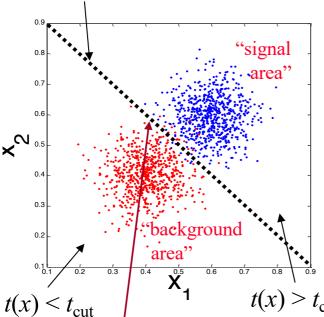
Number of parameters to be determined by MC: M<sup>2</sup> nonlinear test statistic: e.g. neural networks – function in  $\bar{x}$ -space.



#### Linear test statistics



$$t(x) = t_{\text{cut}}$$
 Ansatz:  $t(\overline{x}) = a_0 + \sum_{i=1}^n a_i x_i = a_0 + \overline{a}^T \overline{x}$ 



 $t(\bar{x})$  a Fisher linear discriminant see e.g.

R.A. Fisher, *Annals of Eugenics* **7** (1936) 179; R.A. Fisher: Contributions to mathematical statistics, 1950.

A choice of  $\bar{a}$  gives certain pdf's  $g(t \mid H_0) \& g(t \mid H_1)$ 

 $t(x) > t_{\text{cut}}$  choose  $a_i$ 's to maximize "separation" between

n = 2, a  $t_{\text{cut}}$  can be interpreted  $g(t | H_0) \& g(t | H_1)$  a straight line in the  $x_1 - x_2$  plane.  $\Rightarrow$  must define

"separation" more exactly !!

Data  $\bar{x} = (x_1,...,x_n)$  have mean values & covariance matrix

$$(\mu_k)_i = \int x_i f(\bar{x} | H_k) dx_1 ... dx_n \quad i, j = 1, ..., n$$

$$(V_k)_{ij} = \int (x - \mu_k)_i (x - \mu_k)_j f(\bar{x} \mid H_k) dx_1 ... dx_n \quad k = 0,1$$

Each hypothesis has certain mean & variance of  $t(\bar{x})$ 

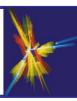
$$\tau_k = \int t(\bar{x})g(t(\bar{x})|H_k)dt = \bar{a}^T \bar{\mu}_k$$

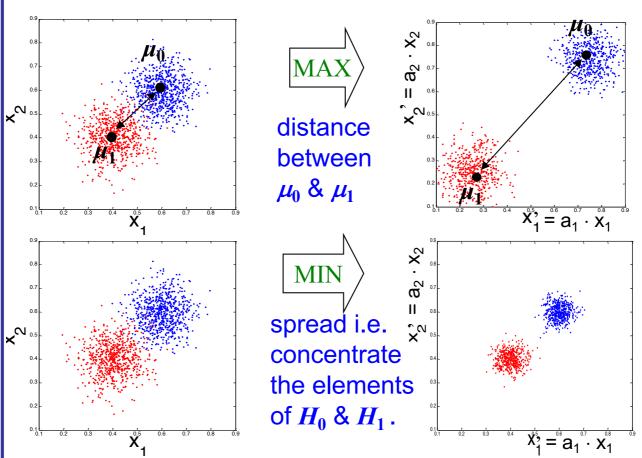
$$\Sigma_k^2 = \int (t(\bar{x}) - \tau_k)^2 g(t(\bar{x}) | H_k) dt = \bar{a}^T \mathbf{V}_k \bar{a}$$

To optimize separation: large  $|\tau_0 - \tau_1|$ , small  $\Sigma_0^2 \& \Sigma_1^2$  (pdf's tightly concentrated about well separated means).



#### Fisher discriminant





Fisher defined as a measure of separation (taking both into account):  $J(\bar{a}) = (\tau_0 - \tau_1)^2 / (\Sigma_0^2 + \Sigma_1^2)$ 

Numerator & denominator can be written as

$$(\tau_0 - \tau_1)^2 = \sum_{i,j=1}^n a_i a_j (\mu_0 - \mu_1)_i (\mu_0 - \mu_1)_j = \sum_{i,j=1}^n a_i a_j B_{ij} = \overline{a}^T \mathbf{B} \, \overline{a}$$
$$\Sigma_0^2 + \Sigma_1^2 = \sum_{i,j=1}^n a_i a_j (\mathbf{V}_0 + \mathbf{V}_1)_{ij} = \overline{a}^T \mathbf{W} \, \overline{a}$$

This gives  $J(\bar{a}) = \frac{\bar{a}^T \mathbf{B} \bar{a}}{\bar{a}^T \mathbf{W} \bar{a}} = \frac{\text{separation between classes}}{\text{sum of variances within classes}}$ 

set 
$$\partial J(\overline{a})/\partial a_i = 0$$
, for all  $i \Rightarrow \overline{a} \propto \mathbf{W}^{-1}(\overline{\mu}_0 - \overline{\mu}_1)$ 

Above Fisher's linear discriminant function & corresponding test statistic a Fisher discriminant.



#### Fisher discriminant



Note coefficients  $a_i$  only determined up to arbitrary scale factor. Use scale  $(t_{cut})$  or & offset  $a_0$  to fix purity and efficiency. Larger  $t_{cut}$  can consimprove signal purity with cate good efficiency loss & vice-versa.

arbitrary 0.8  $t_{cut}$  increases scale  $(t_{cut})$  0.7  $t_{cut}$  and 0.6  $t_{cut}$  can  $cut_{0.5}$   $cut_{0$ 

To determine coefficients  $a_i$ , and need matrix **W** & expectation

values  $\mu_{(0,1)i}$ . Usually estimated from a set of training data e.g. from a MC simulation. Important point being that one doesn't need to determine joint pdf's  $f(\bar{x} | H_0)$  &  $f(\bar{x} | H_1)$  as n-dimensional histograms, only means & covariances. Instead of determining  $M^n$  parameters now only n(n+1)/2.

Maximizing  $J(\overline{a}) = \frac{(\tau_0 - \tau_1)^2}{(\Sigma_0^2 + \Sigma_1^2)}$  with fixed  $\tau_0 \& \tau_1$  same as minimizing  $\Sigma_0^2 + \Sigma_1^2 = E_0[(t(\overline{x}) - \tau_0)^2] + E_1[(t(\overline{x}) - \tau_1)^2]$ 

 $(E_k$  denotes expectation value under hypothesis  $H_k$ )

 $\rightarrow$  maximizing Fisher's  $J(\bar{a})$  kind of **least square** problem. (more about parameter determination using least squares later)

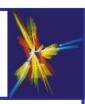
Special Fisher discriminants:  $f(\bar{x}|H_k)$  multidimensional Gaussians with same covariance matrix  $V_0 = V_1 \equiv V$ 

$$f(\bar{x} \mid H_k) = \frac{1}{(2\pi)^{n/2} |\mathbf{V}|^{1/2}} \exp\left[-\frac{1}{2}(\bar{x} - \bar{\mu}_k)^T \mathbf{V}^{-1}(\bar{x} - \bar{\mu}_k)\right], \quad k = 0, 1$$

and the Fisher discriminant:  $t(\bar{x}) = a_0 + (\bar{\mu}_0 - \bar{\mu}_1)^T \mathbf{V}^{-1} \bar{x}$ 



#### Fisher discriminant



Recall likelihood ratio (maximum efficiency at given purity)

$$r = \frac{f(\bar{x} | H_0)}{f(\bar{x} | H_1)} = \exp\left[\frac{1}{2}(\bar{x} - \bar{\mu}_1)^T \mathbf{V}^{-1}(\bar{x} - \bar{\mu}_1) - \frac{1}{2}(\bar{x} - \bar{\mu}_0)^T \mathbf{V}^{-1}(\bar{x} - \bar{\mu}_0)\right]$$

$$r = \exp[(\overline{\mu}_0 - \overline{\mu}_1)^T \mathbf{V}^{-1} \overline{x} - \frac{1}{2} \overline{\mu}_0^T \mathbf{V}^{-1} \overline{\mu}_0 + \frac{1}{2} \overline{\mu}_1^T \mathbf{V}^{-1} \overline{\mu}_1] \propto e^t$$

i.e.  $t(\bar{x}) \propto \ln r + \text{const.}$  t given by monotonic function of r  $\Rightarrow$  Fisher discriminant = likelihood ratio i.e. most optimal.

NB! for non-equal  $V_i$  or non-Gaussian pdfs, no longer true.

Multidimensional Gaussian with equal covariance matrices also gives simple expressions for posterior hypotheses probabilities  $f(\bar{x} \mid H)_{\pi}$ 

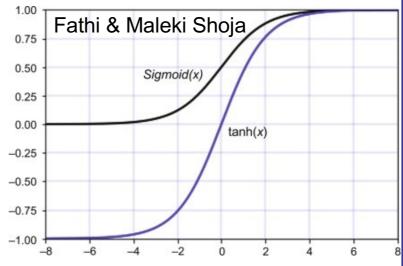
$$\begin{array}{ll} \text{hypotheses probabilities} & f(\overline{x}\,|\,H_0)\pi_0\\ \text{e.g.} & P(H_0\,|\,\overline{x}) = \overline{f(\overline{x}\,|\,H_0)\pi_0 + f(\overline{x}\,|\,H_1)\pi_1} = \overline{1 + (\pi_1/\pi_0 r)}, \\ \text{(Bayes' theorem)} & \overline{f(\overline{x}\,|\,H_0)\pi_0 + f(\overline{x}\,|\,H_1)\pi_1} = \overline{1 + (\pi_1/\pi_0 r)}, \end{array}$$

where  $\pi_0 \& \pi_1$  prior probabilities for  $H_0 \& H_1$ .

Combining this with above expression for r, one gets

$$P(H_0 \mid \bar{x}) = \frac{1}{1 + e^{-t}} \equiv s(t)$$
 the **logistic sigmoid** function (see figure below)

Logistic sigmoid a very common neural network activation function (= giving the weight of the node in the network). Another common one is the hyperbolic tangent:



$$P(H_0|\bar{x}) = \tanh(t) = (e^t - e^{-t})/(e^t + e^{-t})$$



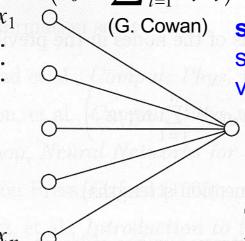
#### **Neural networks**



Optimal decision boundary mostly not "straight line" (due to non-equal  $\mathbf{V}_i$  or non-Gaussian pdfs), use instead nonlinear test statistics: **neural networks**, **boosted decision trees...** 

Artificial Neural Networks (ANN) today used in many fields medical imaging, pattern recognition, financial forecasting etc... but also in physics. Assume  $t(\bar{x})$  to have the form

 $t(\overline{x}) = s\left(a_0 + \sum_{i=1}^n a_i x_i\right) \text{ where } s(\cdot) \text{ activation function (often logistic sigmoid or } \tanh(t))$ 



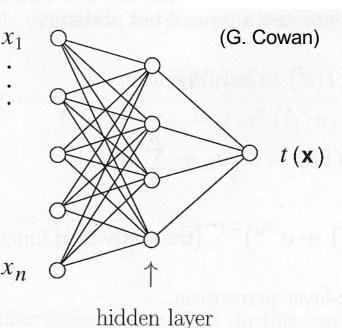
**single-layer perceptron** a test statistic of this form.  $\bar{x}$  input values represented as a set

of nodes.  $s(\cdot)$  monotonic  $\Rightarrow t(\bar{x})$  corresponds to a linear test output node (could statistic.

be more than one)

input layer

Generalized to multilayer:
perceptrons (MLP) with
one or several <u>hidden</u>
layers. Each hidden layer
has m nodes  $(h_1, ..., h_m)$ .
Usually connection are
restricted: value of a given
node only depends on  $x_n$ nodes in previous layer,
a feed-forward network.





#### **Neural networks**



In a two-layer perceptron output defined by

$$t(\overline{x}) = s\left(a_0 + \sum_{i=1}^m a_i h_i(\overline{x})\right) \text{ where } h_i(\overline{x}) = s\left(w_{i0} + \sum_{j=1}^n w_{ij} x_j\right)$$

 $t(\bar{x})$  **non-linear test statistic**, since non-linear function of inputs  $x_i$ .  $a_i$  and  $w_i$  **connection strengths** or **weights** (n+1)m free parameters for n input parameters. With more nodes, ANN closer to optimal but more parameters to fix.

Parameters determined by minimizing an error function

$$\varepsilon = E_0[(t - t^{(0)})^2] + E_1[(t - t^{(1)})^2]$$

where  $t^{(0)}$  &  $t^{(1)}$  target values e.g. 0 & 1 for logistic sigmoid, equivalent to least square method for Fisher discriminants.

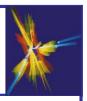
To determine parameters minimizing  $\varepsilon$ , iterative numerical methods, network <u>training</u> or <u>learning</u>, used. In practice expectation values replaced by means computed from training samples, i.e. MC simulations. Learning starts with random initial values for weights and proceed evaluating function using training data. Weights adjusted to minimize  $\varepsilon$  by different methods, i.e. error-back-propagation.

**Theorem:** MLP with a single hidden layer having sufficient number of nodes can approximate arbritrarily well optimal decision boundary; Leshno et al., Neural networks 6 (1993) 861

Advantages/disadvantages of using a single hidden layer with many nodes over many hidden layers not known exactly but seems reasonable to assume that many layers have better performance & more stable than a single layer  $\Rightarrow$  deep learning (networks with 10's or 100's of hidden layers)



#### Nonlinear test statistics



# Nonlinear test statistics Focus here on supervised learning i.e. answer for each training event known

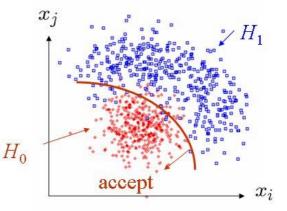
The optimal decision boundary may not be a hyperplane,

 $\rightarrow$  nonlinear test statistic  $t(\vec{x})$ 

Multivariate statistical methods are a Big Industry:

- · artificial neural networks
- boosted decision trees
- generative Al

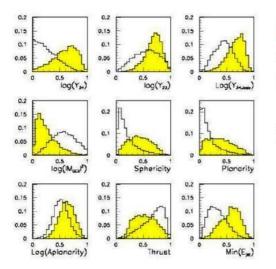
. . .



Physics can benefit from progress in Machine learning

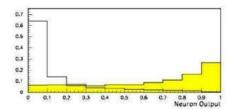
## Neural network example from LEP II

Signal:  $e^+e^- \rightarrow W^+W^-$  (often 4 well separated hadron jets) yellow Background:  $e^+e^- \rightarrow qqgg$  (4 less well separated hadron jets) white



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

Neural network output does better...



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

Machine learning in python: sklearn, https://scikit-learn.org/



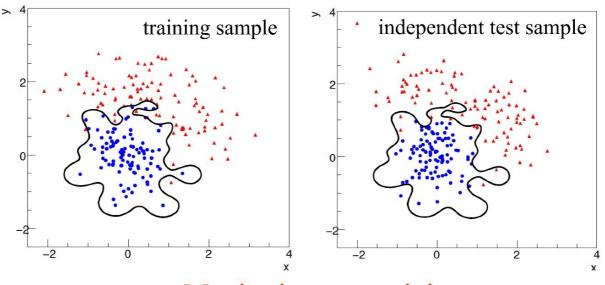
#### **Overtraining**



## Overtraining

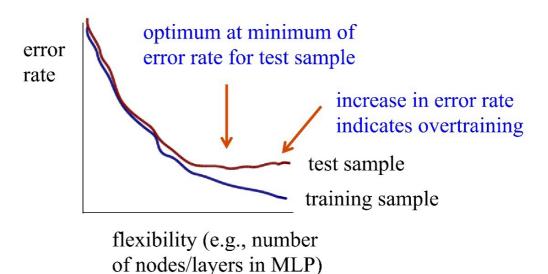
Including more parameters in a classifier makes its decision boundary increasingly flexible, e.g., more nodes/layers for a neural network.

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").



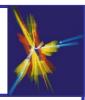
## Monitoring overtraining

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





#### **Boosted decision trees**



## For multiclassification problems can work better in practice than e.g. neural networks

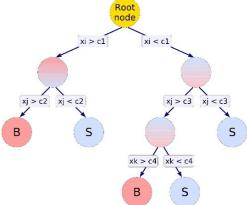
A training sample of signal and background data is repeatedly split by successive cuts on its input variables.

Order in which variables used based on best separation between signal and background.

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

Resulting set of cuts is a 'decision tree'.

Tends to be sensitive to fluctuations in training sample.



## Boosted decision trees

sklearn implementation: Randomforestregressor

Boosting combines a number classifiers into a stronger one; improves stability with respect to fluctuations in input data.

To use with decision trees, increase the weights of misclassified events and reconstruct the tree.

Iterate  $\rightarrow$  forest of trees (perhaps  $\geq$  1000). For the *m*th tree,

application

$$T_m(\vec{x}) = egin{cases} 1 & \vec{x} \text{ in signal acceptance region Jan Welti, PhD} \\ -1 & \text{otherwise} \end{cases}$$
 thesis in physics, HU-P-D235

Define a score  $\alpha_m$  based on error rate of mth tree.

(2017)

Boosted tree = weighted sum of the trees:  $T(\vec{x}) = \sum_{m} \alpha_m T_m(\vec{x})$ 

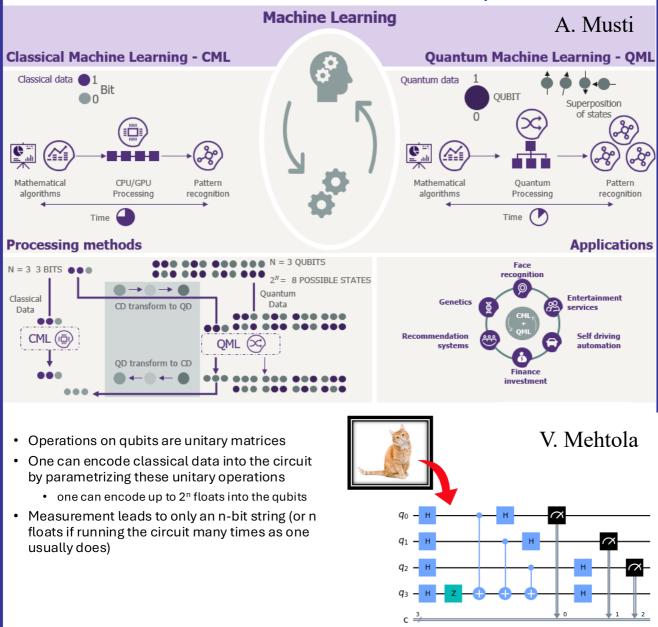
Algorithms: AdaBoost (Freund & Schapire), ε-boost (Friedman).

bagging = draw bootstraped samples from training sample, create tree for each & finally combine them



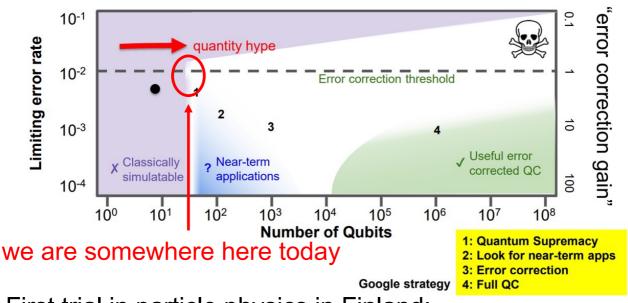


- When data points projected in higher and higher dimensions, becomes harder for a classical computer to deal with it. Even if it does, it takes too much time.
- Sometimes, classical machine learning algorithms too taxing for classical computers.
- Solution: quantum computers? They use superposition & entanglement to solve problems (potentially) much faster than their classical counterparts.





## **Need Both Quality and Quantity**

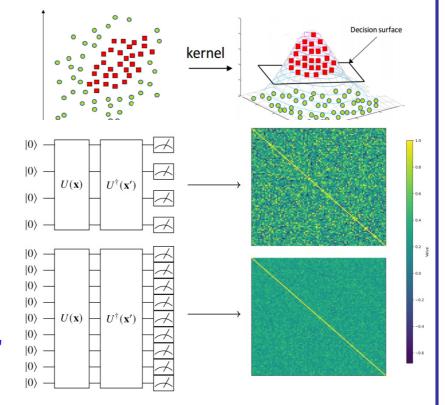


First trial in particle physics in Finland:

V. Mehtola MSc thesis 2025

quantum kernel for selecting Vector Boson Scattering events from background

Need sufficient number of qubits to remove "noise" (= limit error rate) to be competitive with classical machine learning.



**Figure 3:** A figure for visualizing the general phenomenon of concentration (exponential or not). On the left, there are two quantum kernel circuits with same feature map scaled to four and nine qubits, respectively. To their right is their respective Gram matrix. For the nine-qubit version, off-diagonal variance is visibly lower. In the case that this scaling is exponential in the number of qubits, the feature map  $U(\mathbf{x})$  does not scale well at least for the given data and pre-processing.



#### 2D example (G. Cowan)



## 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).

 $f(x_1|x_2) \sim$  Gaussian, different means for s/b, Gaussians have same  $\sigma$ , which depends on  $x_2$ ,  $f(x_2) \sim$  exponential, same for both s and b,  $f(x_1, x_2) = f(x_1|x_2) f(x_2)$ :

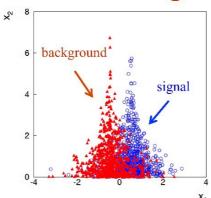
$$f(x_1, x_2|\mathbf{s}) = \frac{1}{\sqrt{2\pi}\sigma(x_2)} e^{-(x_1 - \mu_\mathbf{s})^2/2\sigma^2(x_2)} \frac{1}{\lambda} e^{-x_2/\lambda}$$

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

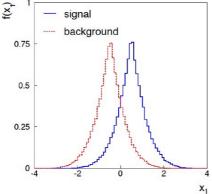
$$\sigma(x_2) = \sigma_0 e^{-x_2/\xi}$$

## Joint and marginal distributions of $x_1, x_2$

0.75







Distribution  $f(x_2)$  same for s, b. So does  $x_2$  help discriminate between the two event types?

signal background

#### 2D example (G. Cowan)



## Likelihood ratio for 2D example

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

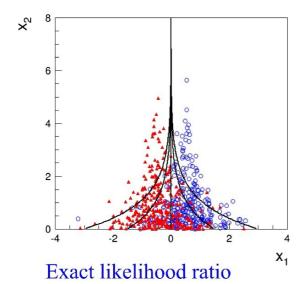
$$t(x_1, x_2) = \frac{f(x_1, x_2|\mathbf{s})}{f(x_1, x_2|\mathbf{b})}$$

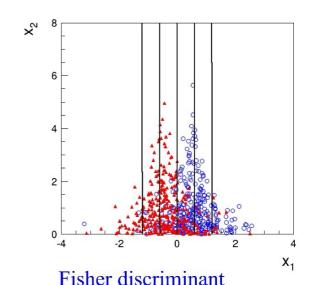
Equivalently we can use any monotonic function of this as a test statistic, e.g.,

$$\ln t = \frac{\frac{1}{2}(\mu_{\rm b}^2 - \mu_{\rm s}^2) + (\mu_{\rm s} - \mu_{\rm b})x_1}{\sigma_0^2 e^{-2x_2/\xi}}$$

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

## Contours of constant MVA output

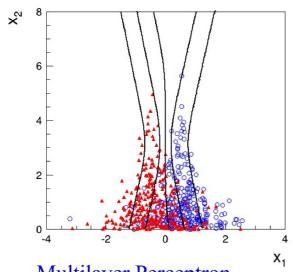




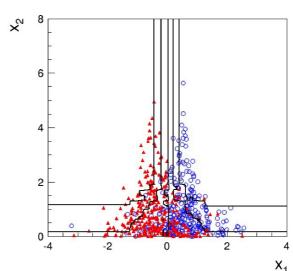
#### 2D example (G. Cowan)



## Contours of constant MVA output



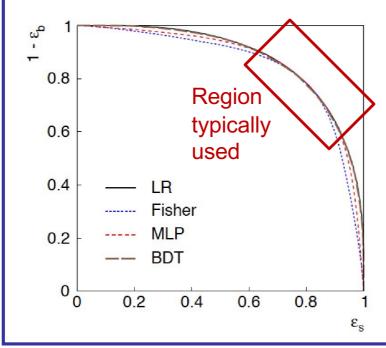
Multilayer Perceptron
1 hidden layer with 2 nodes



Boosted Decision Tree 200 iterations (AdaBoost)

Training samples: 10<sup>5</sup> signal and 10<sup>5</sup> background events

## ROC curve

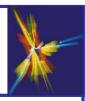


ROC = "receiver operating characteristic" (term from signal processing).

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.





## Multivariate analysis discussion

First clean the input of faulty, problematic or inadequate data For all methods, need to check:

> Sensitivity to statistically unimportant variables (best to drop those that don't provide discrimination);

> Level of smoothness in decision boundary (sensitivity to over-training)

Given the test variable, next step is e.g., select n events and estimate a cross section of signal:  $\hat{\sigma}_s = (n-b)/\varepsilon_s L$ 

Now need to estimate systematic error...

If e.g. training (MC) data \neq Nature, test variable is not optimal, but not necessarily biased.

But our estimates of background b and efficiencies would then be biased if based on MC. (True also for 'simple cuts'.)

But in a cut-based analysis it may be easier to avoid regions where untested features of MC are strongly influencing the decision boundary.

Look at control samples to test joint distributions of inputs.

Try to estimate backgrounds directly from the data (sidebands).

The purpose of the statistical test is often to select objects for further study and then measure their properties.

> Need to avoid input variables that are correlated with the properties of the selected objects that you want to study. (Not always easy; correlations may be poorly known.)

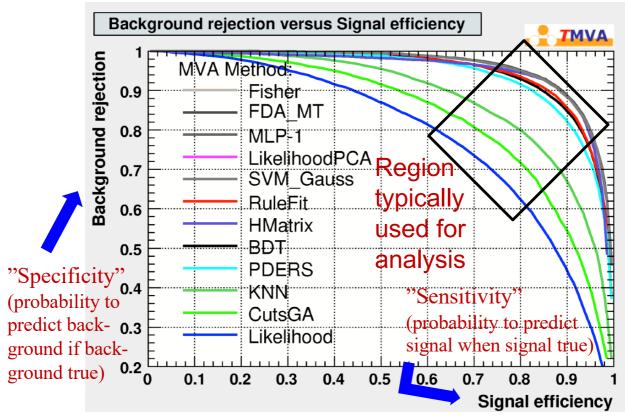
Also decide whether "hard" classification (defining explicitly signal & background regions) is necessary or whether "soft" classification (that gives simply signal & background probabilities) is sufficient.



#### **Multivariate methods (TMVA)**



Toolkit for Multivariate Analysis (TMVA) included in ROOT-package enables usage of many techniques simultaneously:

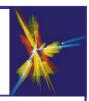


Multivariate/machine learning techniques included:

- Neural networks
- Deep networks
- Multilayer perceptron (MLP)
- Boosted/Bagged decision trees (BDT)
- Support Vector Machine (SVM)
- Fisher and Function discriminant analysis (FDA)
- Multidimensional probability density estimation (PDE) range-search approach (PDERS)
- Multidimensional k-nearest neighbour classifier (KNN)
- Predictive learning via rule ensembles (RuleFit)
- Projective likelihood estimation (PDE approach)
- Rectangular cut optimisation

# HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

#### **Multivariate methods**



# 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*, Springer, 2001

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

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

PDG review on machine learning:

https://pdg.lbl.gov/2025/reviews/rpp2024-rev-machine-learning.pdf

## 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

https://root.cern.ch/

## **SKLEARN** - Machine Learning in Python

- Simple and efficient tools for predictive data analysis
- Built on NumPy, SciPy, and matplotlib https://scikit-learn.org/stable/

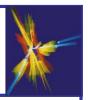
**KERAS** – Python interface to many different neural network packages <a href="https://keras.io/">https://keras.io/</a>

**PyTorch** – Python interface for neural networks in Torch programme package <a href="https://pytorch.org/">https://pytorch.org/</a>

Dedicated course at University of Helsinki: DATA11002 Introduction to Machine Learning 5 ECTS (period 2/2025)

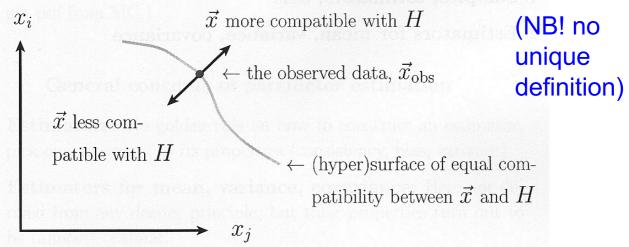


#### Goodness-of-fit tests



Want to substanciate how well hypothesis H compatible with observed data (without reference to an alternative).

Hypothesis H predict  $f(\bar{x}|H)$  for data vector  $\bar{x} = (x_1, ... x_n)$ Observe a single point in  $\bar{x}$ -space:  $\bar{x}_{obs}$ . What can be said about validity of H in light of data?  $\to$  decide what parts of  $\bar{x}$ -space less compatible with H than observed point  $\bar{x}_{obs}$ .



Construct test statistic  $t(\bar{x})$  whose value reflects level of compatibility between  $\bar{x}$  & H. e.g. high  $t \rightarrow$  data less compatible with H low  $t \rightarrow$  data more compatible with H.

## Express **goodness-of-fit** by giving **P-value**:

P = probability to observe data  $\bar{x}$  (or  $t(\bar{x})$ ) having equal or lesser compatibility with H than  $\bar{x}_{obs}$  (or  $t(\bar{x}_{obs})$ ).

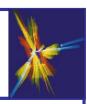
NB! *P*-value not be interpreted as probability that *H* true !! *P*-value also <u>observed significance level/confidence level</u>

If H true, then (for continuous  $\bar{x}$ ) P uniform in [0,1] If H not true, then pdf of P (usually) peaked close to 0.

Does small P-value really mean H is false? No, P-value = probability to obtain such a result "by chance" if H is true.



#### Goodness-of-fit example



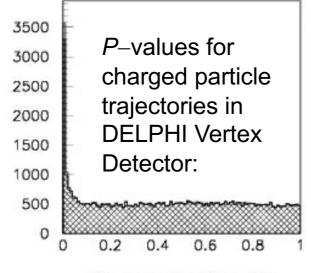
Example: fit of charged particle trajectory in magnetic field Primary aim, obtain parameters of trajectory ("helix"). Goodness-of-fit variables to determine compatibility of trajectory segments/hits to charged particle trajectory but also for checks of covariance matrix for helix parameters. Much simplified, if trajectory consists of N hits measured in  $R\phi$  & M hits measured in Z, then track fit  $\approx$  to minimize:

$$\chi^{2} = \sum_{i=1}^{N} (R\phi_{i} - f_{R\phi}(R_{i}))^{2} / \sigma(R\phi_{i})^{2} \& \chi^{2} = \sum_{i=1}^{M} (z_{i} - f_{z}(R_{i}))^{2} / \sigma(z_{i})^{2}$$

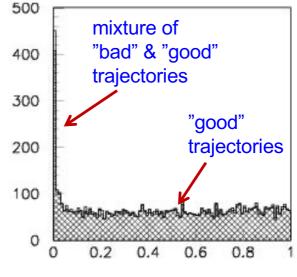
where  $f_{R\phi/z}(R)$  describes dependence w.r.t. to radius, R. Typical for particle physics.  $N_{dof} = N-3$  (solenoidal field) & M-2. Here multiple scattering & energy loss ignored. Calculate corresponding P-value from observed  $\chi^2$  value:

 $P = \int_{\chi^2}^{\infty} f(z; N_{dof}) dz$  where  $f(z; N_{dof})$  chi-square pdf for  $N_{dof}$  degrees of freedom.

K.Österberg, PhD thesis, HIP-1998-01



VD track probability in RPhi

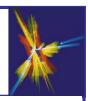


VD track probability in Rz

Distribution should be flat: raise towards 0 indicate that there is a class of hits with a slightly worse resolution & peak at 0 misassociated hits or badly reconstructed tracks.



#### HELSINGFORS UNIVERSITET Kolmogorov-Smirnov test



# To make goodness-of-fit tests of distributions, the **Kolmogorov–Smirnov (K-S) test** is often used.

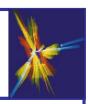
- By construction, it was developed for the comparison of functions, but is widely also used for binned data.
- The K-S test is predominantly sensitive to the shape of the distribution, much more than the  $\chi^2$ -test. It's also normalization independent, which is not the  $\chi^2$ -test. Used mainly in three ways:
- test whether data sample follows certain distribution
- test whether two data sample could be consistent to come from the same distribution
- test whether data in two histograms are consistent with each other
  - Kolmogorov distance is defined as:
    - $dist_{K} = Max|F_{parent}(x) F_{test}(x)|$ , where F is cumulative distribution for parent and test normalized to one.
    - distance is also know as Kolmogorov test statistic.
  - Probability that dist<sub>K</sub> ≥ X is given by Kolmogorov distribution function:

$$P(dist_k \ge x) = 2\sum_{i=1}^{\infty} (-1)^{i-1} e^{-2i^2x^2}$$

NB! true if n large (> ~10), where n is size of data sample.



#### Kolmogorov-Smirnov test

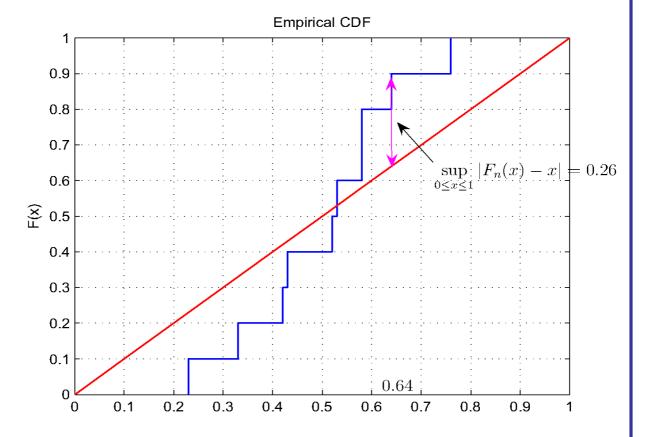


Example: 10 random numbers between 0 and 1: 0.58, 0.42, 0.52, 0.33, 0.43, 0.23, 0.58, 0.76, 0.53, 0.64 Are they uniformly distributed in [0,1]?

i.e. 
$$H_0$$
:  $F_{\text{test}}(x) = F_0(x) = x$ 

$$\max_{0 \le x \le 1} |F_{data}(x) - F_{test}(x)| = \max_{1 \le i \le n} \left( |F_n(X_i^-) - x|, |F_n(X_i^+) - x| \right)$$

cumulative distribution function value before/after ith jump

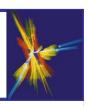


KS probability:  $dist_k \rightarrow dist_k \cdot \sqrt{n}$  so here  $0.26 \cdot \sqrt{10} \approx 0.82$  hypothesis testing at significance level  $\alpha = 0.05$ :

$$1 - P(c \ge x) = 0.05 \Rightarrow c = 1.35 \ (0.82 \Rightarrow P\text{-value} = 17 \ \%)$$

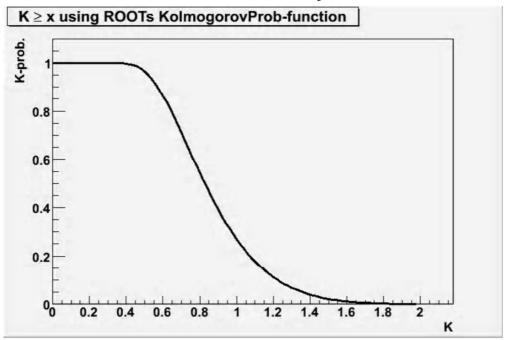
Here  $dist_k \cdot \sqrt{n} = 0.82 < c = 1.35 \Rightarrow$  hence hypothesis random numbers uniform in [0,1] can not be rejected.

#### Kolmogorov-Smirnov test



## KS-test

Kolmogorov distribution is slow to integrate and there are plenty of algorithms to calculate Kolmogorovs probability: eg. ROOTs TMath library, matlab, octave, mathematica, java...



When comparing how two one-dimensional distibutions differ, Kolmogorov probability is same when calculated by replacing:

$$dist_K \rightarrow dist_K \sqrt{\frac{nn'}{n+n'}}$$

where n and n' are the number of entries in parent 2 and test distributions.

testing if two

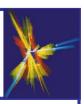
data sample

originate

from same

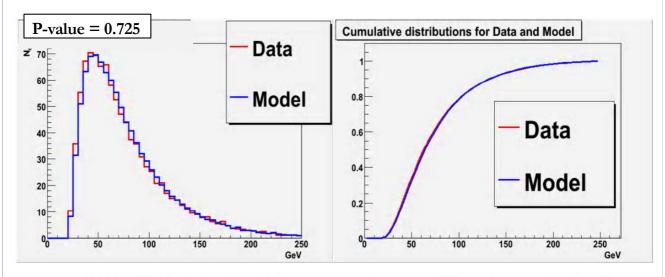
distribution!!

#### **Kolmogorov-Smirnov test**

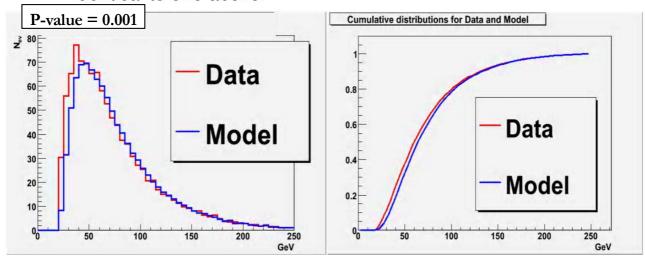


## **KS-test**

 P-value is calculated comparing how often simulated experiment has larger K<sub>dist</sub> than K-distance is in data vs model.

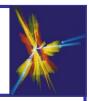


- WARNING: Few bad bins can make your KS-value really bad!
   Check your distributions, understand your measurements! One bin with 4-sigma deviation can make your KS close to zero.
- Eg. Low energy region badly modelled. Otherwise distribution identical to one above.



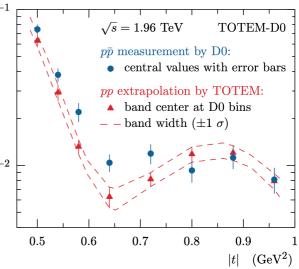


### Kolmogorov-Smirnov test

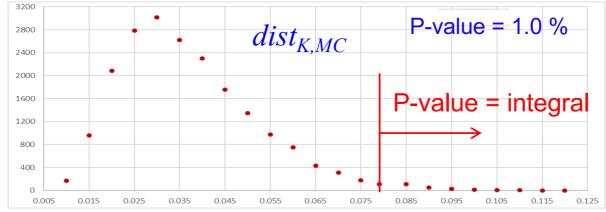


An additional usage of Kolmogorov-Smirnov is to test the compatibility of two sets of measurements with uncertainties using MC methods. In this case one tests whether the shape of the distributions of the two data sets are compatible with each other.

What is the probability that data set 1 (blue, taking its uncertainties into account) with uncertainties of data set 2 (red) would give  $dist_K >= dist_{K,data}^{10}$  i.e. the  $dist_K$  between data set 1 & 2? Answer: 1.0 %

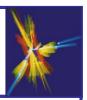


- 1. determine  $dist_{K data}$  between data 1 and 2 (= 0.079)
- 2. generate data set 1' (1") by adding to the central value of each data set 1 point, a random gaussian according to point uncertainty in data set 1 (2)
- 3. calculate  $dist_{KMC}$  between data set 1' and 1"
- 4. Repeat step 2 and 3 to get a distribution of  $dist_{KMC}$
- 5. Calculate probability that  $dist_{K,MC} >= dist_{K,data}$





#### Pearson's χ<sup>2</sup> test



Let's examine another goodness-of-fit test applicable to distribution of a variable x. Assume histogram of observed x-values in N bins. Suppose number of entries in bin i is  $n_i$  & number of expected entries  $v_i$ , then the test statistic

$$\chi^2 = \sum_{i=1}^{N} (n_i - v_i)^2 / v_i$$

reflects level of agreement between observed & expected histogram. Above test based on **Pearson's**  $\chi^2$  **statistic.** If data  $\bar{n} = (n_1, ..., n_N)$  consists of N independent Poisson variables with mean values  $\bar{v} = (v_1, ..., v_N)$  & all  $v_i$  not too small (rule of thumb: all  $v_i \ge 5$ ) then test statistic  $\chi^2$  follow chi-square pdf for N degrees of freedom. Holds regardless of distribution for variable x ( $\chi^2$  test distribution free).

Standard deviation of Poisson distribution  $\sqrt{\nu_i}$  so  $\chi^2$  test sum of difference squared between observed & expected values, measured in units of standard deviations squared.

Corresponding P-value given by the observed  $\chi^2$  using  $P = \int_{\chi^2}^{\infty} f(z; N) dz$  where f(z; N) chi-square pdf for N degrees of freedom.

(google "chi square calculator" to get cumulative  $\chi^2$  function) Recall that for chi-square pdf, expectation value E[z] = N $\rightarrow$  often give  $\chi^2/N$  as a measure of level of agreement. However better to give  $\chi^2$  & N separately...

$$\chi^2$$
 = 15,  $N$  = 10  $\rightarrow$   $P$ -value = 13 %  $\chi^2$  = 150,  $N$  = 100  $\rightarrow$   $P$ -value = 0.09 %

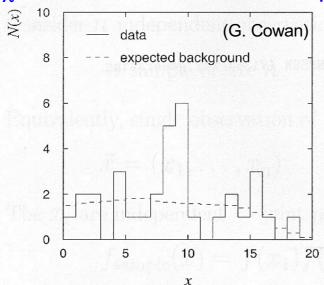
If  $n_{\text{tot}} = \sum_{i=1}^{N} n_i$  fixed then  $n_i$ 's multinomial with  $v_i = p_i n_{\text{tot}}$  &  $\chi^2 = \sum_{i=1}^{N} (n_i - p_i n_{\text{tot}})^2 / p_i n_{\text{tot}}$  follows chi-square pdf with N-1 dof (if all  $p_i n_{\text{tot}} \gg 1$ ).



#### Pearson's χ² test



Even if amount of data too small for requirement all  $v_i \ge 5$  to be fulfilled, one can still construct  $\chi^2$  statistic as long as all  $v_i$ 's > 0. Will no longer follow chi-square distribution &  $\chi^2$  statistic distribution will depend on pdf of variable x. e.g.

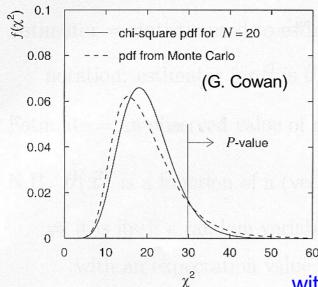


histogram (to the left) with peak gives  $\chi^2$  = 29.8 for N = 20 dof.

but ...all bins have  $\nu_i$  < 5. here  $\chi^2$  statistic will not follow a normal chisquare distribution.

Pearson's  $\chi^2$  still used as test statistic, but must generate corresponding  $f(\chi^2)$  using MC methods to get P-value.

- generate  $n_i$  from Poissonian with mean  $v_i$ , for i = 1, ..., N,
- compute & record  $\chi^2$  value into a histogram,
- repeat steps above many times to get  $f(\chi^2)$ ,
- compute *P*-value as probability of  $\chi^2 \ge \chi^2_{\text{obs}}$  using  $f(\chi^2)$ .



using  $f(\chi^2)$  from MC:

*P*-value = 11 %

using chi-square pdf: *P*-value = 7.3 %

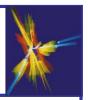
The  $\chi^2$  test not very sensitive to presence of peaks.

NB! binning ambiguity of  $\chi^2$  test.

P-value can vary significantly with binning for small data samples.



### Significance of a signal



A simple type of goodness-of-fit test carried out to judge whether a discrepancy between data & expectation is sufficiently significant to merit a claim for a new discovery.

One observes *n* events; these can then consist of:

 $n_b$  events from a known process (background)

 $n_s$  events from a possible new process (signal)

If  $n_b$  &  $n_s$  Poisson random variables with means  $v_b$  &  $v_s \Rightarrow n = n_s + n_b$  also Poissonian with mean  $v = v_b + v_s$ :

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

Suppose  $v_b$  = 0.5 & one observes  $n_{obs}$  = 5.

Should one claim evidence for a new discovery? hypothesis  $H: v_s = 0$ , i.e. only background events present.

$$P$$
 - value =  $P(n \ge n_{\text{obs}}) = \sum_{n=n_{\text{obs}}}^{\infty} P(n; v_s = 0, v_b) =$ 

 $1 - \sum_{n=0}^{n_{\rm obs}-1} \frac{v_b^n}{n!} e^{-v_b} = 1.7 \cdot 10^{-4} \ (\neq P(v_s = 0)!) \ \text{for } v_s = 0 \ \text{hypothesis}$ 

Typical misunderstandings:

a misleading (but often occuring) estimate...

estimate for mean  $v: n_{obs} = 5$ 

estimated standard deviation of  $n: \sqrt{n} = 2.2 \implies$ 

"measured signal" =  $n_{\rm obs} - \nu_b$  = 4.5  $\pm$  2.2 i.e.  $\sim$  2 $\sigma$  from 0.

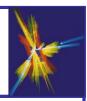
Wanted: probability for Poisson variable with mean  $v_b$  = 0.5 to give 5 or more (answer: 0.017 %)

Misunderstanding implied: probability for variable with  $\mu$  = 4.5 &  $\sigma$  = 2.2 to give 0 or less (Gaussian answer: 2.2 %)

NB! difference disappear when  $v_b$  »1, i.e. Poisson ~ Gaussian.



#### Significance of a signal

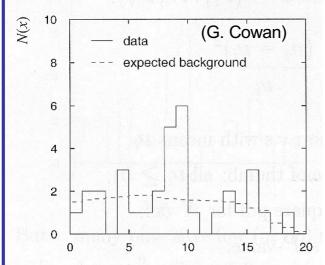


another pitfall: believe  $v_b$  to be God-given; in reality must study influence of systematic uncertainties on  $v_b$ . if e.g.  $\sigma(v_b \operatorname{syst.}) = 0.3 \Rightarrow P(n \ge 5; v_b = 0.8, v_s = 0) = 0.14 \%$   $\Rightarrow$  conservative approach: report P-value corresponding to  $\pm 1 \sigma$  in disfavourable direction of result, here P = 0.14 %

When can one claim observation? Often require P-value corresponding to  $5\sigma(3\sigma)$  of a Gaussian distribution for new ("known") phenomena i.e.  $\leq 2.9 \cdot 10^{-7}$  ( $\leq 1.35 \cdot 10^{-3}$ ). Particle physics convention:  $>5\sigma$  = discovery,  $>3\sigma$  = evidence

### Significance of a peak:

In addition to counting events, one measures also x for each event; data viewed as a histogram as function of x:



histogram of observed data & theoretical expectation. Each bin a Poisson variable, since each x value independent (or multinomial if the expectation is normalised to the total number of observed events).

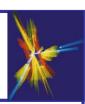
Is the peak significant?

In the two bins of peak, 11 observed events,  $v_b = 3.2$ ,

$$P(n \ge 11; \ v_b = 3.2, \ v_s = 0) = 0.05 \%$$

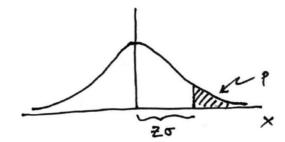
but ... did one know where to look for the peak?  $\rightarrow$  give  $P(n \ge 11)$  in any 2 adjacent bins ("look elsewhere" effect). Is observed width consistent with expected x resolution? How many bins  $\times$  distributions did one examine? Did one adjust selections to "enhance" the peak? What about the side-bands of the peak, are they too low?





## 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.



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

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

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

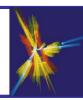
Rejecting  $H_0$  with Z > 5 constitutes "discovering" new physics.

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:

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

### Likelihood ratio of experiment (G. Cowan)



## 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

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

and for the "signal plus background" hypothesis  $(H_1)$  it is

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

where  $\pi_s$  and  $\pi_b$  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

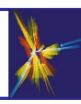
$$Q = -2\ln\frac{L_{s+b}}{L_b} = -s + \sum_{i=1}^{n} \ln\left(1 + \frac{s}{b} \frac{f(\mathbf{x}_i|\mathbf{s})}{f(\mathbf{x}_i|\mathbf{b})}\right)$$

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.

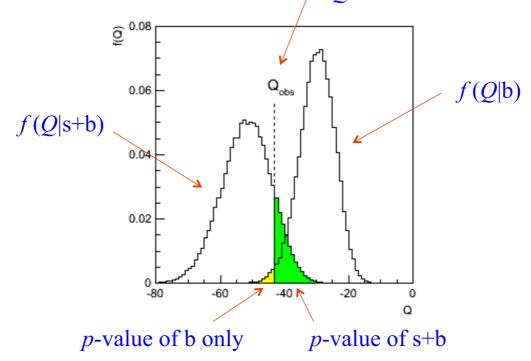
### HELSINGFORS UNIVERSITET Likelihood ratio of experiment (G. Cowan)



## Distribution of Q

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

Suppose in real experiment *O* is observed here.



# Systematic uncertainties

Up to now we assumed all parameters were known exactly.

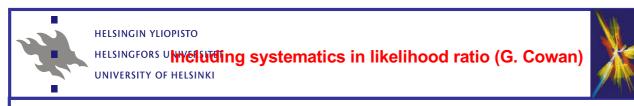
In practice they have some (systematic) uncertainty.

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  $\sigma_b$  is the estimated uncertainty.



# 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.

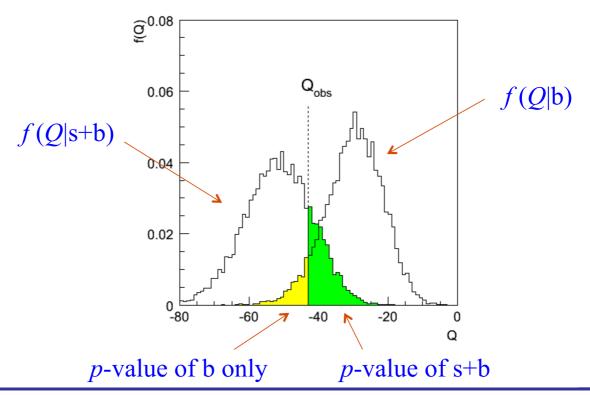
But we can obtain the Bayesian model average:

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

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  $Q_{\rm obs}$ .

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





#### **ANOVA**



What if number of data samples to compare is large?

Can test each sample pair with t– (basic method to test if two samples have same mean or not),  $\chi^2$  or KS-test... but easily get into trouble since probability of (at least) one pair to differ significantly naturally gets high (like "peak" search in large # of histograms, "peaks" are found).

Analysis of Variance (ANOVA): collection of statistical methods to treat comparison of 3 or more data samples. Originally developed by R. Fisher in 1920's and 30's.

Comparison can be made on one or several variables (one-way ANOVA and two(three)-way or factorial ANOVA)

Typical assumptions: independent samples, normal distributed samples, equality of variances...

Logic of ANOVA:

$$V_{tot}[x] = V_{stat}[x] + V_{diff\ sample}[x]$$

ndf of  $\chi^2$  distribution describing sum of variances

$$\operatorname{ndf}_{tot}[x] = \operatorname{ndf}_{stat}[x] + \operatorname{ndf}_{diff\ sample}[x]$$

$$F = \frac{\text{variance of sample means}}{\text{mean of within - sample variances}}$$

variable F follows F-distribution. Cumulative F-distribution to get P-value. If replace "mean of within-sample variances" with "expected variation of sample means" then F follows  $\chi^2$  distribution with ndf = # of samples -1



#### **ANOVA**



**Example 11.1** A barmaid at the Bull's Head is working to pay her way through university. To get a little intellectual exercise while pulling pints, she keeps track of the drinking habits of the various UMIST departments whose staff drink in the pub. Here are her findings from the previous week:

Physics	Maths	Chemistry
21	6	8
16	10	6
17	13	4
18	13	5
	1	9

The data are usually arranged in columns, as in the table above:

Total number of staff:

number of staff: Column 1 Column 2 ... Column 
$$r$$

$$N = \sum_{j=1}^{r} n_j.$$
  $x_{1,1}$   $x_{1,2}$  ...  $x_{1,r}$  
$$x_{2,1}$$
  $x_{2,2}$  ...  $x_{2,r}$  
$$\vdots$$
 
$$\vdots$$

Then one does the following.

$$X_{n1,1}$$
  $X_{n2,2}$   $X_{nr,r}$ 

a) Compute the mean of each column separately:

$$m_j = \frac{\sum_{k=1}^{n_j} x_{k,j}}{n_j}$$
. "column" = sample

b) Compute the mean of all the entries in the table

$$m = \frac{\sum_{j=1}^{r} \left[\sum_{k=1}^{n_j} x_{k,j}\right]}{N}$$
$$= \frac{\sum_{j=1}^{r} n_j m_j}{N}.$$

c) Compute a measure of the variation between the columns. It is called  $SS_b$ , the "sum of squares between columns", and defined by

$$SS_b = \sum_{j=1}^{r} n_j (m_j - m)^2.$$

d) Compute a measure of the variation within the various columns. It is called  $SS_w$ , the "sum of squares within the columns", and is given by

$$SS_w = \sum_{j=1}^r \left[ \sum_{k=1}^{n_j} (x_{k,j} - m_j)^2 \right]$$



#### **ANOVA**



It's messy, but not difficult, to show that the two sums-of-squares account for all the variation in the whole table:

$$SS_t = SS_b + SS_w = \sum_{j=1}^r \left[ \sum_{k=1}^{n_j} (x_{k,j} - m)^2 \right]$$

Here the subscript 't' stands for "total";  $SS_t$  is the sum-of-squared-deviations from the mean for the entire table.

One then assembles these results into a standard table that looks like:

			Average sum of squares
Type of Variation	Sum of Squares	deg. of freedom	(estimated variances)
Between Samples	$SS_b$	(r-1)	$SS_b/(r-1)$
Within Samples	$SS_w$	$\sum_{j=1}^{r} (n_j - 1)$ $= (N - r)$	
		=(N-r)	$SS_w/(N-r)$
Total	$SS_t$	(N-1)	_

Finally, one does a test to see whether the mean variation between columns,  $SS_b/(r-1)$  is significantly bigger than the mean variation within columns  $SS_w/(N-r)$ . In the usual way one computes a statistic,

$$F = \frac{SS_b/(r-1)}{SS_w/(N-r)}$$
, Find P-value: google "F-distribution applet"

then looks the value up (in special tables) to see if it is large enough to reject the null hypothesis. These tables depend on the number of degrees of freedom in the sums, so one consults the table for F((r-1), (N-r)).

### Example of beer drinking:

$$N = 14$$
,  $r = 3$ 

			Average sum of squares
Type of Variation	Sum of Squares	deg. of freedom	$(estimated\ variances)$
Between Samples	386.9	2	193.5
Within Samples	102	11	8.5
Total	488.9	13	

 $\Rightarrow$  P-value = 1.2 · 10<sup>-4</sup>

and the test stat. is  $F \approx 22.8$ . The attached table shows that if F(2, 12) exceeds 3.89 we can reject the null with 95% confidence, so we can safely conclude that the members the various groups do not drink the same amount in the Bull's Head.

Scheffe test statistic: expected maximal statistical difference (MSD) of two samples (if the underlying distribution is the same):

MSD = 
$$\sqrt{(r-1) * \overline{V[x]} * F(P = 0.05) * (\frac{1}{n_1} + \frac{1}{n_2})}$$