

# Statistical Methods in Particle Physics

## Day 1: Introduction



清华大学高能物理研究中心  
2010年4月12—16日



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

# Outline of lectures

## Day #1: Introduction

Review of probability and Monte Carlo

Review of statistics: parameter estimation

## Day #2: Multivariate methods (I)

Event selection as a statistical test

Cut-based, linear discriminant, neural networks

## Day #3: Multivariate methods (II)

More multivariate classifiers: BDT, SVM, ...

## Day #4: Significance tests for discovery and limits

Including systematics using profile likelihood

## Day #5: Bayesian methods

Bayesian parameter estimation and model selection

# Day #1: outline

## Probability and its role in data analysis

Definition, interpretation of probability

Bayes' theorem

## Random variables and their properties

## A catalogue of distributions

## The Monte Carlo method

## Parameter estimation

Method of maximum likelihood

Method of least squares

## Some statistics books, papers, etc.

G. Cowan, *Statistical Data Analysis*, Clarendon, Oxford, 1998

see also [www.pp.rhul.ac.uk/~cowan/sda](http://www.pp.rhul.ac.uk/~cowan/sda)

R.J. Barlow, *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences*, Wiley, 1989

see also [hepwww.ph.man.ac.uk/~roger/book.html](http://hepwww.ph.man.ac.uk/~roger/book.html)

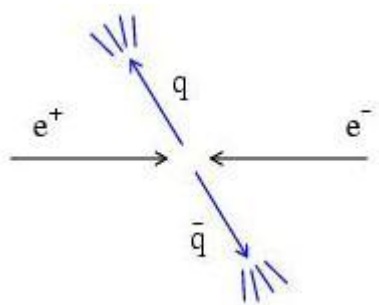
L. Lyons, *Statistics for Nuclear and Particle Physics*, CUP, 1986

F. James., *Statistical and Computational Methods in Experimental Physics*, 2nd ed., World Scientific, 2006

S. Brandt, *Statistical and Computational Methods in Data Analysis*, Springer, New York, 1998 (with program library on CD)

C. Amsler et al. (Particle Data Group), *Review of Particle Physics*, Physics Letters B667 (2008) 1; see also [pdg.lbl.gov](http://pdg.lbl.gov) sections on probability statistics, Monte Carlo

# Data analysis in particle physics



Observe events of a certain type

Measure characteristics of each event (particle momenta, number of muons, energy of jets,...)

Theories (e.g. SM) predict distributions of these properties up to free parameters, e.g.,  $\alpha$ ,  $G_F$ ,  $M_Z$ ,  $\alpha_s$ ,  $m_H$ , ...

Some tasks of data analysis:

Estimate (measure) the parameters;

Quantify the uncertainty of the parameter estimates;

Test the extent to which the predictions of a theory are in agreement with the data ( $\rightarrow$  presence of New Physics?)

# A definition of probability

Consider a set  $S$  with subsets  $A, B, \dots$

For all  $A \subset S, P(A) \geq 0$

$$P(S) = 1$$

If  $A \cap B = \emptyset, P(A \cup B) = P(A) + P(B)$



**Kolmogorov  
axioms (1933)**

Also define conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

# Interpretation of probability

## I. Relative frequency

$A, B, \dots$  are outcomes of a repeatable experiment

$$P(A) = \lim_{n \rightarrow \infty} \frac{\text{times outcome is } A}{n}$$

cf. quantum mechanics, particle scattering, radioactive decay...

## II. Subjective probability

$A, B, \dots$  are hypotheses (statements that are true or false)

$$P(A) = \text{degree of belief that } A \text{ is true}$$

- Both interpretations consistent with Kolmogorov axioms.
- In particle physics frequency interpretation often most useful, but subjective probability can provide more natural treatment of non-repeatable phenomena:

systematic uncertainties, probability that Higgs boson exists,...

# Bayes' theorem

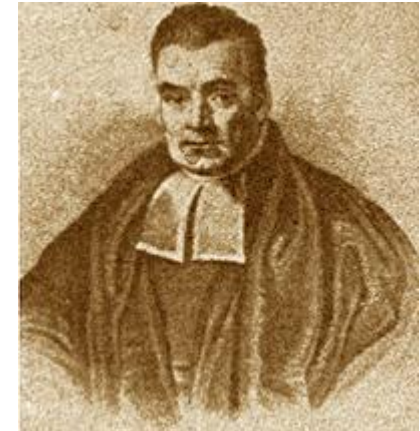
From the definition of conditional probability we have

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{and} \quad P(B|A) = \frac{P(B \cap A)}{P(A)}$$

but  $P(A \cap B) = P(B \cap A)$ , so

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Bayes' theorem



First published (posthumously) by the Reverend Thomas Bayes (1702–1761)

*An essay towards solving a problem in the doctrine of chances*, Philos. Trans. R. Soc. **53** (1763) 370; reprinted in Biometrika, **45** (1958) 293.

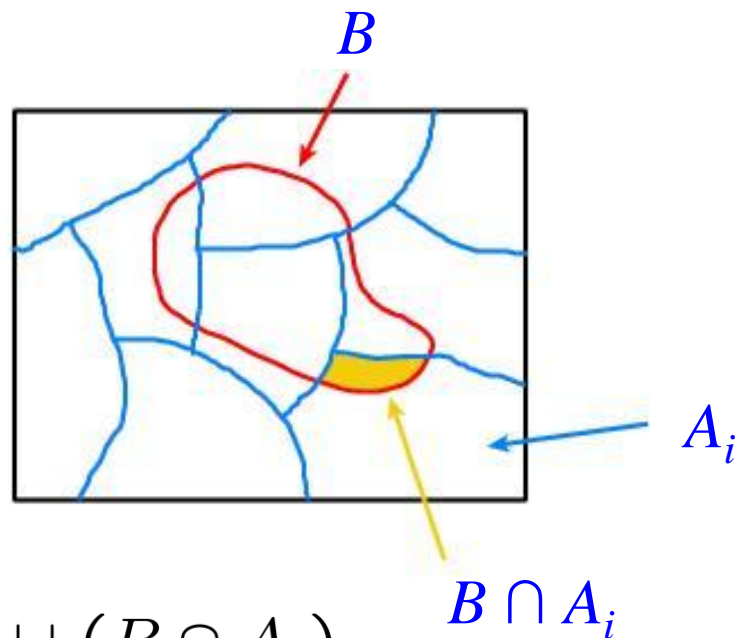


# The law of total probability

Consider a subset  $B$  of the sample space  $S$ ,

divided into disjoint subsets  $A_i$  such that  $\cup_i A_i = S$ ,

$S$



$$\rightarrow B = B \cap S = B \cap (\cup_i A_i) = \cup_i (B \cap A_i),$$

$$\rightarrow P(B) = P(\cup_i (B \cap A_i)) = \sum_i P(B \cap A_i)$$

$$\rightarrow P(B) = \sum_i P(B|A_i)P(A_i) \quad \text{law of total probability}$$

Bayes' theorem becomes

$$P(A|B) = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)}$$

# Random variables and probability density functions

A random variable is a numerical characteristic assigned to an element of the sample space; can be discrete or continuous.

Suppose outcome of experiment is continuous value  $x$

$$P(x \text{ found in } [x, x + dx]) = f(x) dx$$

→  $f(x)$  = probability density function (pdf)

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad x \text{ must be somewhere}$$

Or for discrete outcome  $x_i$  with e.g.  $i = 1, 2, \dots$  we have

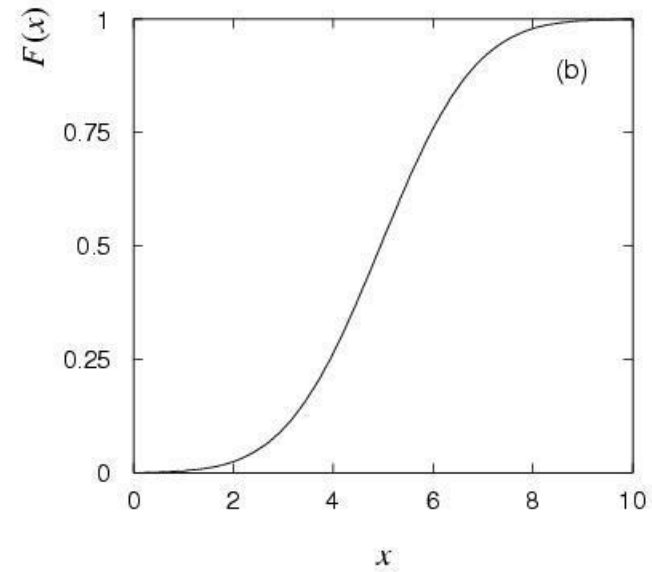
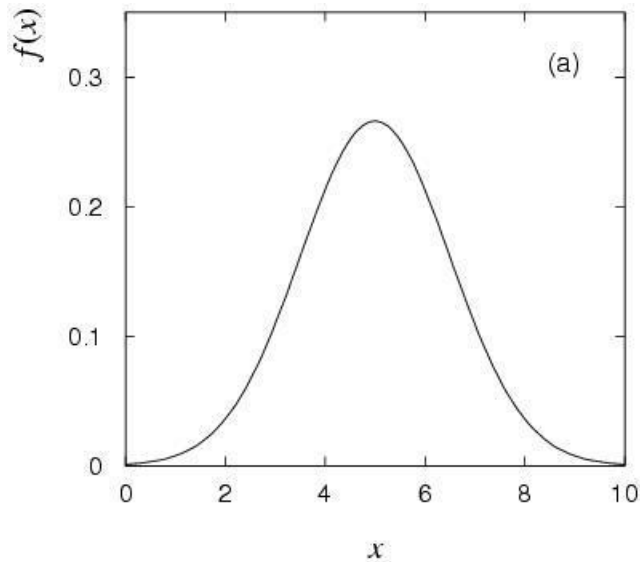
$$P(x_i) = p_i \quad \text{probability mass function}$$

$$\sum_i P(x_i) = 1 \quad x \text{ must take on one of its possible values}$$

# Cumulative distribution function

Probability to have outcome less than or equal to  $x$  is

$$\int_{-\infty}^x f(x') dx' \equiv F(x) \quad \text{cumulative distribution function}$$



Alternatively define pdf with  $f(x) = \frac{\partial F(x)}{\partial x}$

# Other types of probability densities

Outcome of experiment characterized by several values,  
e.g. an  $n$ -component vector,  $(x_1, \dots, x_n)$

→ joint pdf  $f(x_1, \dots, x_n)$

Sometimes we want only pdf of some (or one) of the components

→ marginal pdf  $f_1(x_1) = \int \cdots \int f(x_1, \dots, x_n) dx_2 \cdots dx_n$

$x_1, x_2$  independent if  $f(x_1, x_2) = f_1(x_1)f_2(x_2)$

Sometimes we want to consider some components as constant

→ conditional pdf  $g(x_1|x_2) = \frac{f(x_1, x_2)}{f_2(x_2)}$

# Expectation values

Consider continuous r.v.  $x$  with pdf  $f(x)$ .

Define expectation (mean) value as  $E[x] = \int x f(x) dx$

Notation (often):  $E[x] = \mu \sim$  “centre of gravity” of pdf.

For a function  $y(x)$  with pdf  $g(y)$ ,

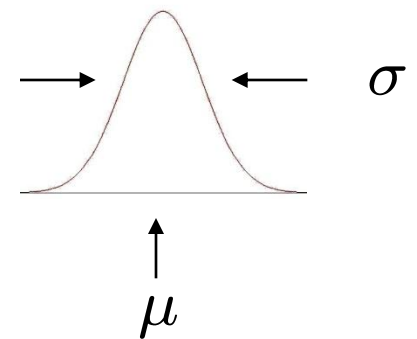
$$E[y] = \int y g(y) dy = \int y(x) f(x) dx \quad (\text{equivalent})$$

Variance:  $V[x] = E[x^2] - \mu^2 = E[(x - \mu)^2]$

Notation:  $V[x] = \sigma^2$

Standard deviation:  $\sigma = \sqrt{\sigma^2}$

$\sigma \sim$  width of pdf, same units as  $x$ .



# Covariance and correlation

Define covariance  $\text{cov}[x,y]$  (also use matrix notation  $V_{xy}$ ) as

$$\text{COV}[x, y] = E[xy] - \mu_x\mu_y = E[(x - \mu_x)(y - \mu_y)]$$

Correlation coefficient (dimensionless) defined as

$$\rho_{xy} = \frac{\text{COV}[x, y]}{\sigma_x\sigma_y}$$

If  $x, y$ , independent, i.e.,  $f(x, y) = f_x(x)f_y(y)$ , then

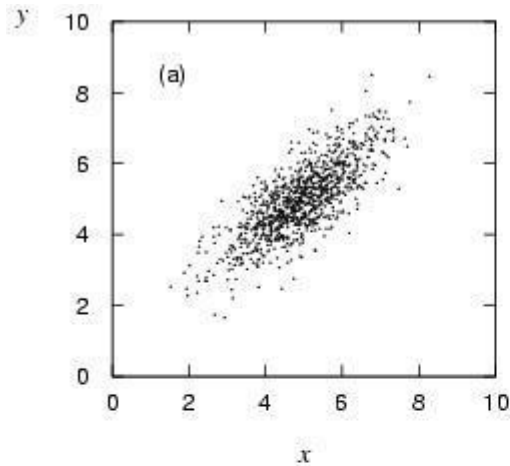
$$E[xy] = \int \int xy f(x, y) dx dy = \mu_x\mu_y$$

→  $\text{COV}[x, y] = 0$       $x$  and  $y$ , ‘uncorrelated’

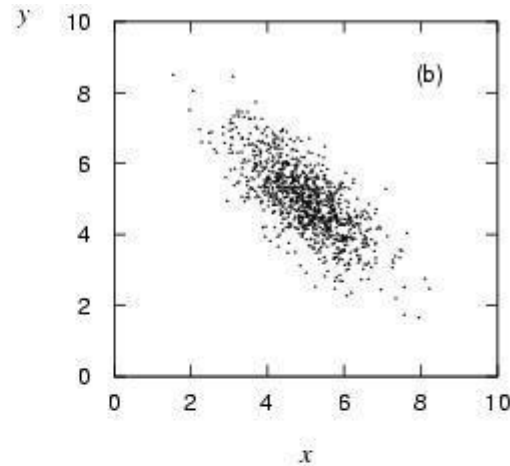
N.B. converse not always true.

# Correlation (cont.)

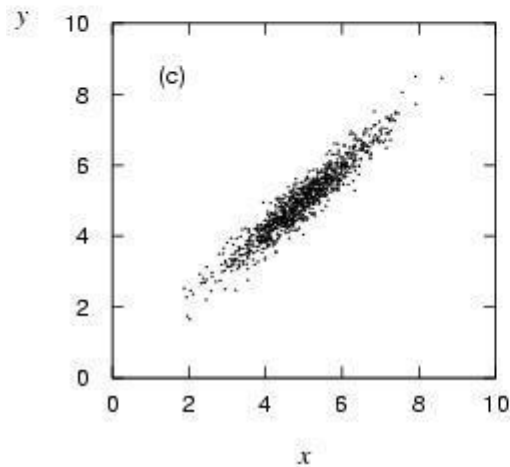
$$\rho = 0.75$$



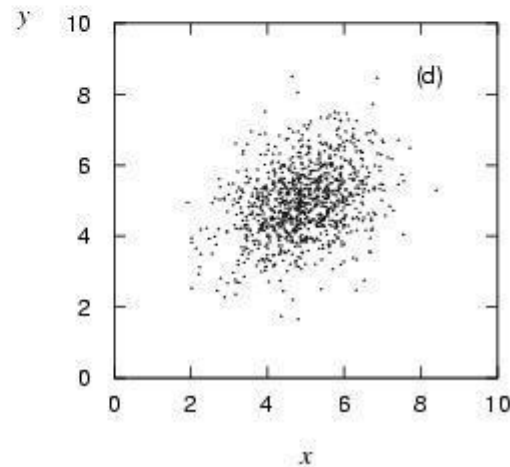
$$\rho = -0.75$$



$$\rho = 0.95$$



$$\rho = 0.25$$



# Some distributions

<u>Distribution/pdf</u>	<u>Example use in HEP</u>
Binomial	Branching ratio
Multinomial	Histogram with fixed $N$
Poisson	Number of events found
Uniform	Monte Carlo method
Exponential	Decay time
Gaussian	Measurement error
Chi-square	Goodness-of-fit
Cauchy	Mass of resonance
Landau	Ionization energy loss



# Binomial distribution

Consider  $N$  independent experiments (Bernoulli trials):

outcome of each is ‘success’ or ‘failure’,  
probability of success on any given trial is  $p$ .

Define discrete r.v.  $n =$  number of successes ( $0 \leq n \leq N$ ).

Probability of a specific outcome (in order), e.g. ‘ssfsf’ is

$$pp(1-p)p(1-p) = p^n(1-p)^{N-n}$$

But order not important; there are  $\frac{N!}{n!(N-n)!}$

ways (permutations) to get  $n$  successes in  $N$  trials, total probability for  $n$  is sum of probabilities for each permutation.

## Binomial distribution (2)

The binomial distribution is therefore

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

random variable          parameters

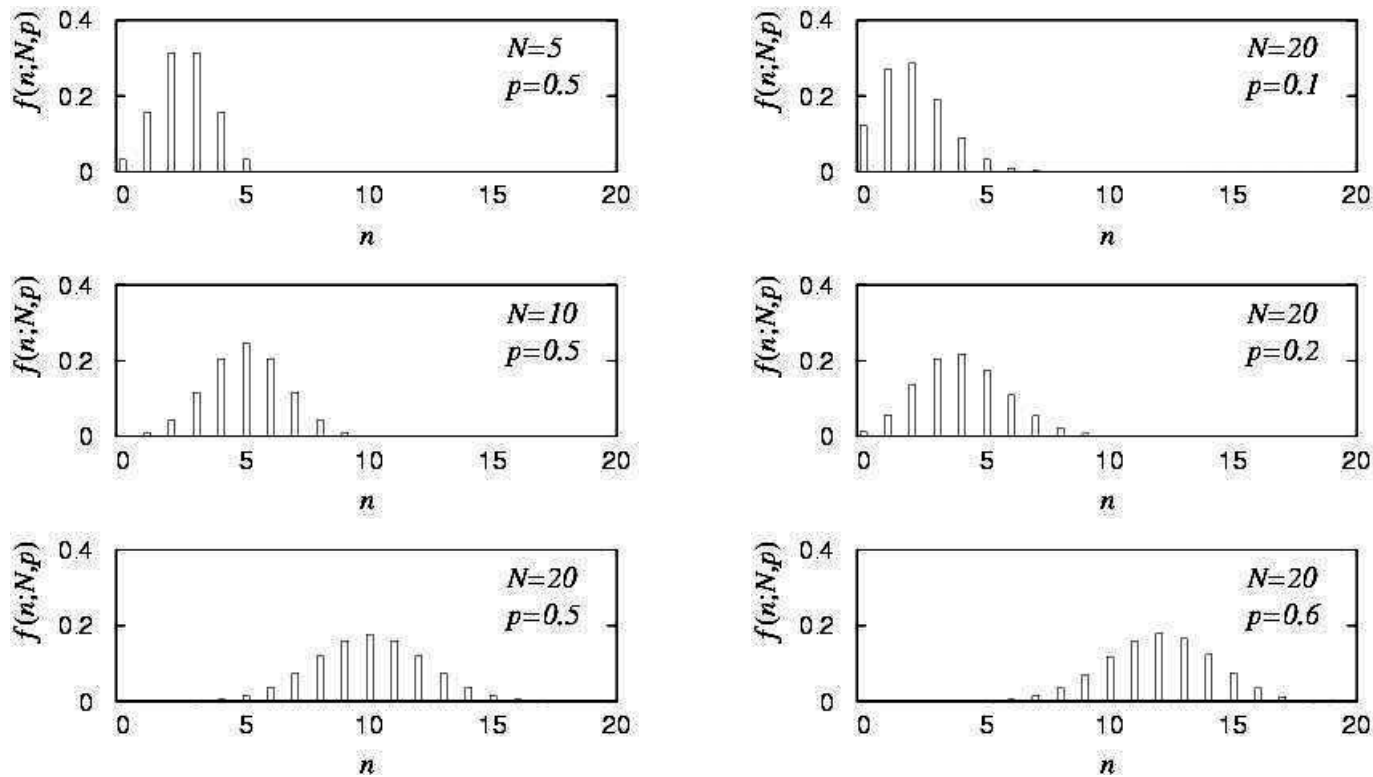
For the expectation value and variance we find:

$$E[n] = \sum_{n=0}^N n f(n; N, p) = Np$$

$$V[n] = E[n^2] - (E[n])^2 = Np(1-p)$$

# Binomial distribution (3)

Binomial distribution for several values of the parameters:



Example: observe  $N$  decays of  $W^\pm$ , the number  $n$  of which are  $W \rightarrow \mu\nu$  is a binomial r.v.,  $p =$  branching ratio.

# Multinomial distribution

Like binomial but now  $m$  outcomes instead of two, probabilities are

$$\vec{p} = (p_1, \dots, p_m), \quad \text{with} \quad \sum_{i=1}^m p_i = 1 .$$

For  $N$  trials we want the probability to obtain:

$n_1$  of outcome 1,  
 $n_2$  of outcome 2,  
...  
 $n_m$  of outcome  $m$ .

This is the multinomial distribution for  $\vec{n} = (n_1, \dots, n_m)$

$$f(\vec{n}; N, \vec{p}) = \frac{N!}{n_1! n_2! \dots n_m!} p_1^{n_1} p_2^{n_2} \dots p_m^{n_m}$$

## Multinomial distribution (2)

Now consider outcome  $i$  as ‘success’, all others as ‘failure’.

→ all  $n_i$  individually binomial with parameters  $N, p_i$

$$E[n_i] = Np_i, \quad V[n_i] = Np_i(1 - p_i) \quad \text{for all } i$$

One can also find the covariance to be

$$V_{ij} = Np_i(\delta_{ij} - p_j)$$

Example:  $\vec{n} = (n_1, \dots, n_m)$  represents a histogram with  $m$  bins,  $N$  total entries, all entries independent.

# Poisson distribution

Consider binomial  $n$  in the limit

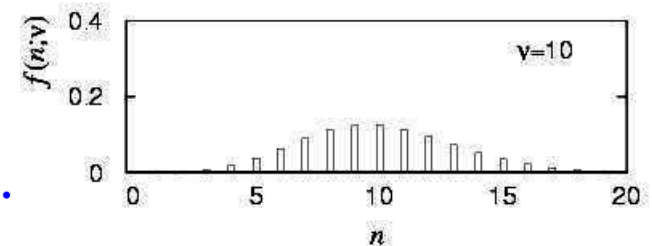
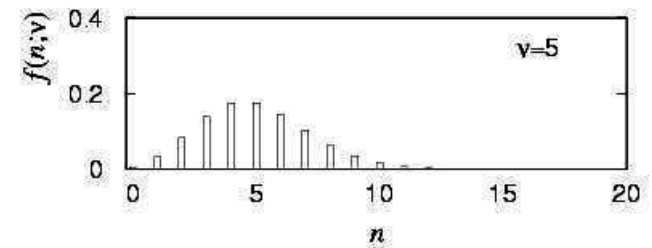
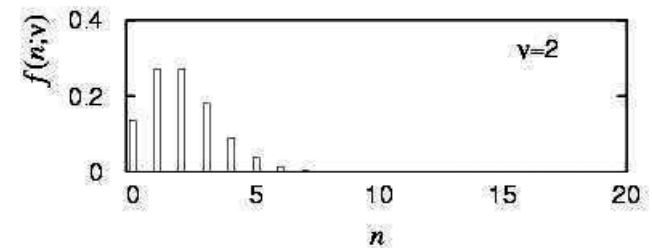
$$N \rightarrow \infty, \quad p \rightarrow 0, \quad E[n] = Np \rightarrow \nu .$$

→  $n$  follows the Poisson distribution:

$$f(n; \nu) = \frac{\nu^n}{n!} e^{-\nu} \quad (n \geq 0)$$

$$E[n] = \nu, \quad V[n] = \nu .$$

Example: number of scattering events  $n$  with cross section  $\sigma$  found for a fixed integrated luminosity, with  $\nu = \sigma \int L dt$ .



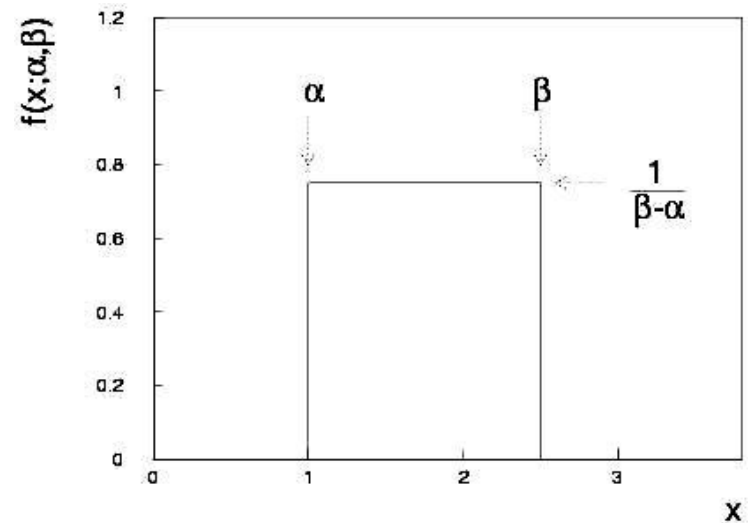
# Uniform distribution

Consider a continuous r.v.  $x$  with  $-\infty < x < \infty$ . Uniform pdf is:

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\beta - \alpha} & \alpha \leq x \leq \beta \\ 0 & \text{otherwise} \end{cases}$$

$$E[x] = \frac{1}{2}(\alpha + \beta)$$

$$V[x] = \frac{1}{12}(\beta - \alpha)^2$$



N.B. For any r.v.  $x$  with cumulative distribution  $F(x)$ ,  $y = F(x)$  is uniform in  $[0,1]$ .

Example: for  $\pi^0 \rightarrow \gamma\gamma$ ,  $E_\gamma$  is uniform in  $[E_{\min}, E_{\max}]$ , with

$$E_{\min} = \frac{1}{2}E_\pi(1 - \beta), \quad E_{\max} = \frac{1}{2}E_\pi(1 + \beta)$$

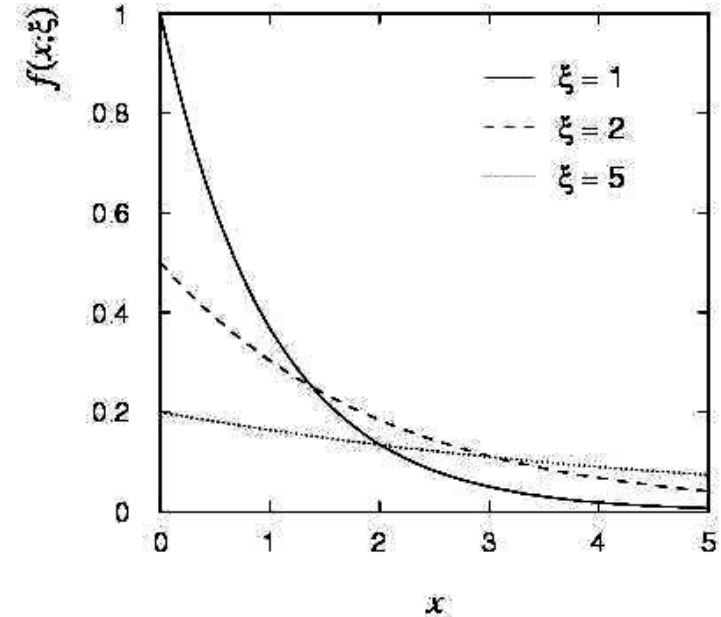
# Exponential distribution

The exponential pdf for the continuous r.v.  $x$  is defined by:

$$f(x; \xi) = \begin{cases} \frac{1}{\xi} e^{-x/\xi} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E[x] = \xi$$

$$V[x] = \xi^2$$



Example: proper decay time  $t$  of an unstable particle

$$f(t; \tau) = \frac{1}{\tau} e^{-t/\tau} \quad (\tau = \text{mean lifetime})$$

Lack of memory (unique to exponential):  $f(t - t_0 | t \geq t_0) = f(t)$



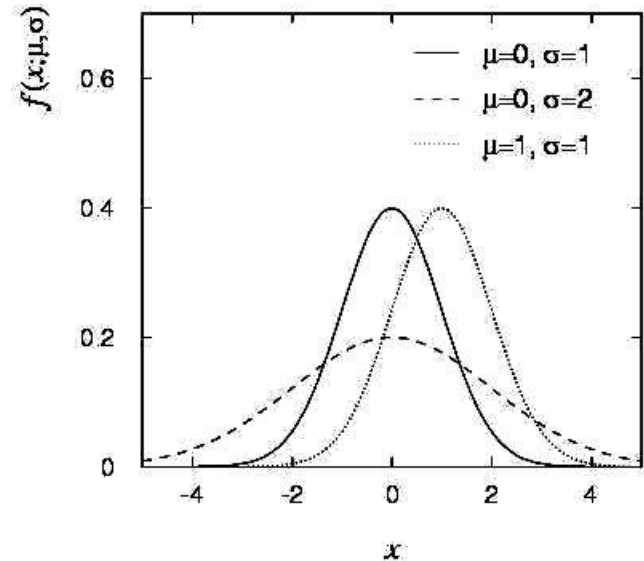
# Gaussian distribution

The Gaussian (normal) pdf for a continuous r.v.  $x$  is defined by:

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$E[x] = \mu \quad (\text{N.B. often } \mu, \sigma^2 \text{ denote mean, variance of any}$$

$$V[x] = \sigma^2 \quad \text{r.v., not only Gaussian.})$$



Special case:  $\mu = 0, \sigma^2 = 1$  ('standard Gaussian'):

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad \Phi(x) = \int_{-\infty}^x \varphi(x') dx'$$

If  $y \sim$  Gaussian with  $\mu, \sigma^2$ , then  $x = (y - \mu) / \sigma$  follows  $\varphi(x)$ .

# Gaussian pdf and the Central Limit Theorem

The Gaussian pdf is so useful because almost any random variable that is a sum of a large number of small contributions follows it. This follows from the Central Limit Theorem:

For  $n$  independent r.v.s  $x_i$  with finite variances  $\sigma_i^2$ , otherwise arbitrary pdfs, consider the sum

$$y = \sum_{i=1}^n x_i$$

In the limit  $n \rightarrow \infty$ ,  $y$  is a Gaussian r.v. with

$$E[y] = \sum_{i=1}^n \mu_i \quad V[y] = \sum_{i=1}^n \sigma_i^2$$

Measurement errors are often the sum of many contributions, so frequently measured values can be treated as Gaussian r.v.s.

## Central Limit Theorem (2)

The CLT can be proved using characteristic functions (Fourier transforms), see, e.g., SDA Chapter 10.

For finite  $n$ , the theorem is approximately valid to the extent that the fluctuation of the sum is not dominated by one (or few) terms.



Beware of measurement errors with non-Gaussian tails.

Good example: velocity component  $v_x$  of air molecules.

OK example: total deflection due to multiple Coulomb scattering. (Rare large angle deflections give non-Gaussian tail.)

Bad example: energy loss of charged particle traversing thin gas layer. (Rare collisions make up large fraction of energy loss, cf. Landau pdf.)

# Multivariate Gaussian distribution

Multivariate Gaussian pdf for the vector  $\vec{x} = (x_1, \dots, x_n)$  :

$$f(\vec{x}; \vec{\mu}, V) = \frac{1}{(2\pi)^{n/2} |V|^{1/2}} \exp \left[ -\frac{1}{2} (\vec{x} - \vec{\mu})^T V^{-1} (\vec{x} - \vec{\mu}) \right]$$

$\vec{x}$ ,  $\vec{\mu}$  are column vectors,  $\vec{x}^T$ ,  $\vec{\mu}^T$  are transpose (row) vectors,

$$E[x_i] = \mu_i, \quad \text{COV}[x_i, x_j] = V_{ij} .$$

For  $n = 2$  this is

$$f(x_1, x_2; \mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \\ \times \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[ \left( \frac{x_1 - \mu_1}{\sigma_1} \right)^2 + \left( \frac{x_2 - \mu_2}{\sigma_2} \right)^2 - 2\rho \left( \frac{x_1 - \mu_1}{\sigma_1} \right) \left( \frac{x_2 - \mu_2}{\sigma_2} \right) \right] \right\}$$

where  $\rho = \text{cov}[x_1, x_2]/(\sigma_1 \sigma_2)$  is the correlation coefficient.

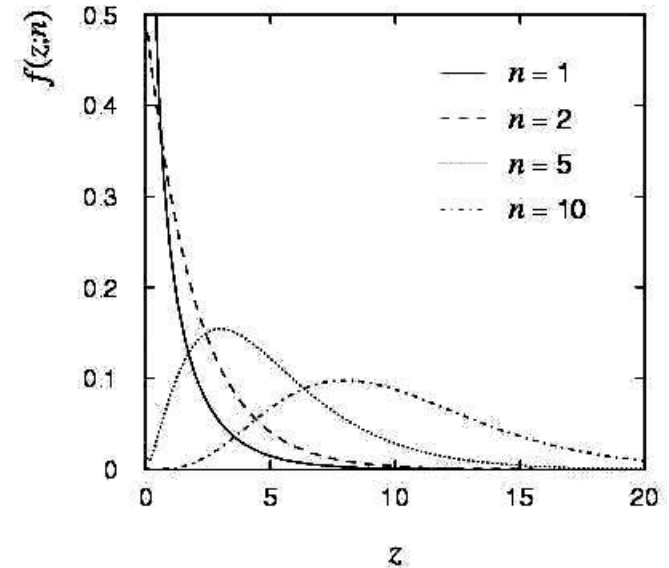
# Chi-square ( $\chi^2$ ) distribution

The chi-square pdf for the continuous r.v.  $z$  ( $z \geq 0$ ) is defined by

$$f(z; n) = \frac{1}{2^{n/2} \Gamma(n/2)} z^{n/2-1} e^{-z/2}$$

$n = 1, 2, \dots$  = number of 'degrees of freedom' (dof)

$$E[z] = n, \quad V[z] = 2n.$$



For independent Gaussian  $x_i$ ,  $i = 1, \dots, n$ , means  $\mu_i$ , variances  $\sigma_i^2$ ,

$$z = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2} \quad \text{follows } \chi^2 \text{ pdf with } n \text{ dof.}$$

Example: goodness-of-fit test variable especially in conjunction with method of least squares.

# Cauchy (Breit-Wigner) distribution

The Breit-Wigner pdf for the continuous r.v.  $x$  is defined by

$$f(x; \Gamma, x_0) = \frac{1}{\pi} \frac{\Gamma/2}{\Gamma^2/4 + (x - x_0)^2}$$

( $\Gamma = 2$ ,  $x_0 = 0$  is the Cauchy pdf.)

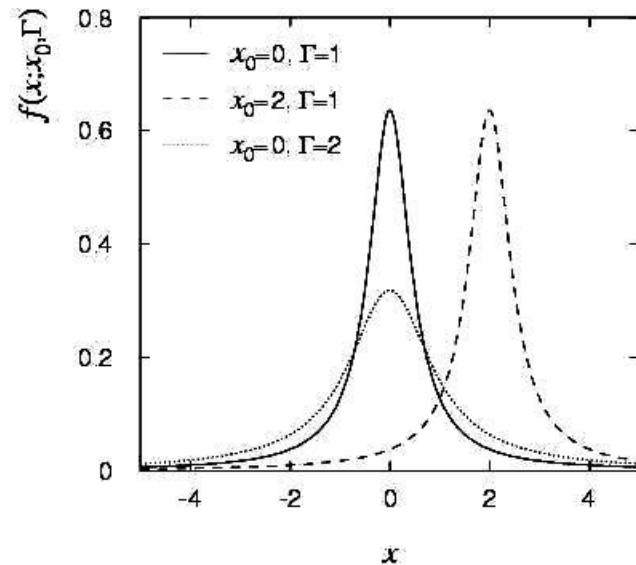
$E[x]$  not well defined,  $V[x] \rightarrow \infty$ .

$x_0 = \text{mode}$  (most probable value)

$\Gamma = \text{full width at half maximum}$

Example: mass of resonance particle, e.g.  $\rho$ ,  $K^*$ ,  $\phi^0$ , ...

$\Gamma = \text{decay rate}$  (inverse of mean lifetime)



# Landau distribution

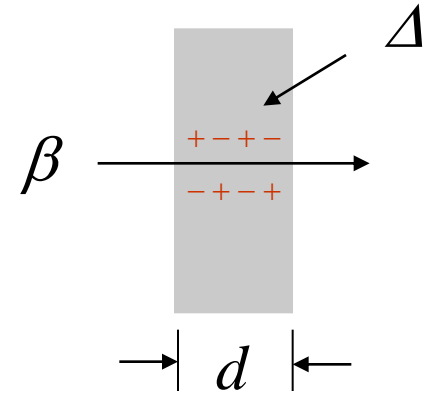
For a charged particle with  $\beta = v/c$  traversing a layer of matter of thickness  $d$ , the energy loss  $\Delta$  follows the Landau pdf:

$$f(\Delta; \beta) = \frac{1}{\xi} \phi(\lambda) ,$$

$$\phi(\lambda) = \frac{1}{\pi} \int_0^\infty \exp(-u \ln u - \lambda u) \sin \pi u \, du ,$$

$$\lambda = \frac{1}{\xi} \left[ \Delta - \xi \left( \ln \frac{\xi}{\epsilon'} + 1 - \gamma_E \right) \right] ,$$

$$\xi = \frac{2\pi N_A e^4 z^2 \rho \sum Z}{m_e c^2 \sum A} \frac{d}{\beta^2} , \quad \epsilon' = \frac{I^2 \exp \beta^2}{2m_e c^2 \beta^2 \gamma^2} .$$



L. Landau, J. Phys. USSR **8** (1944) 201; see also

W. Allison and J. Cobb, Ann. Rev. Nucl. Part. Sci. **30** (1980) 253.

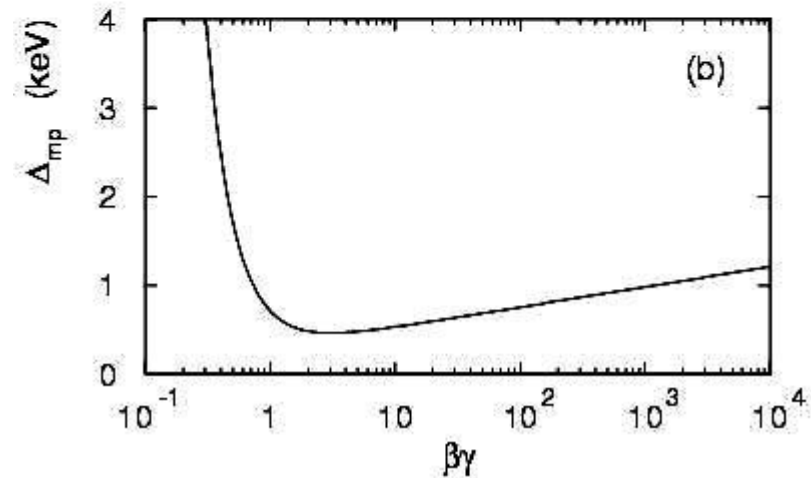
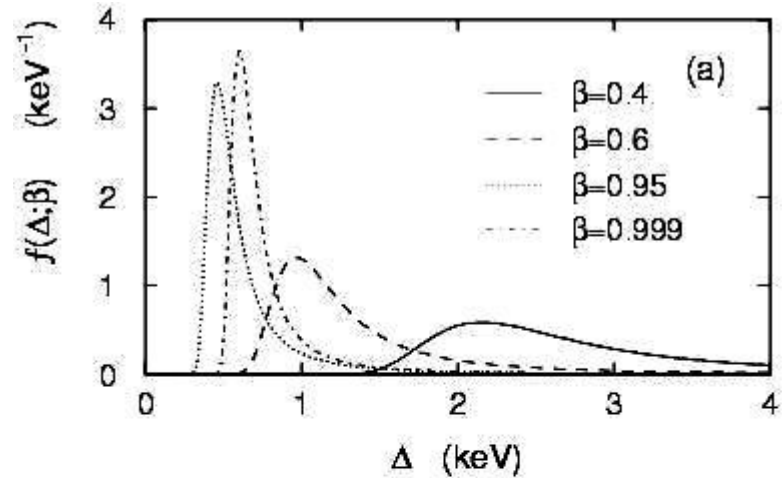
# Landau distribution (2)

Long 'Landau tail'

→ all moments  $\infty$

Mode (most probable value) sensitive to  $\beta$ ,

→ particle i.d.





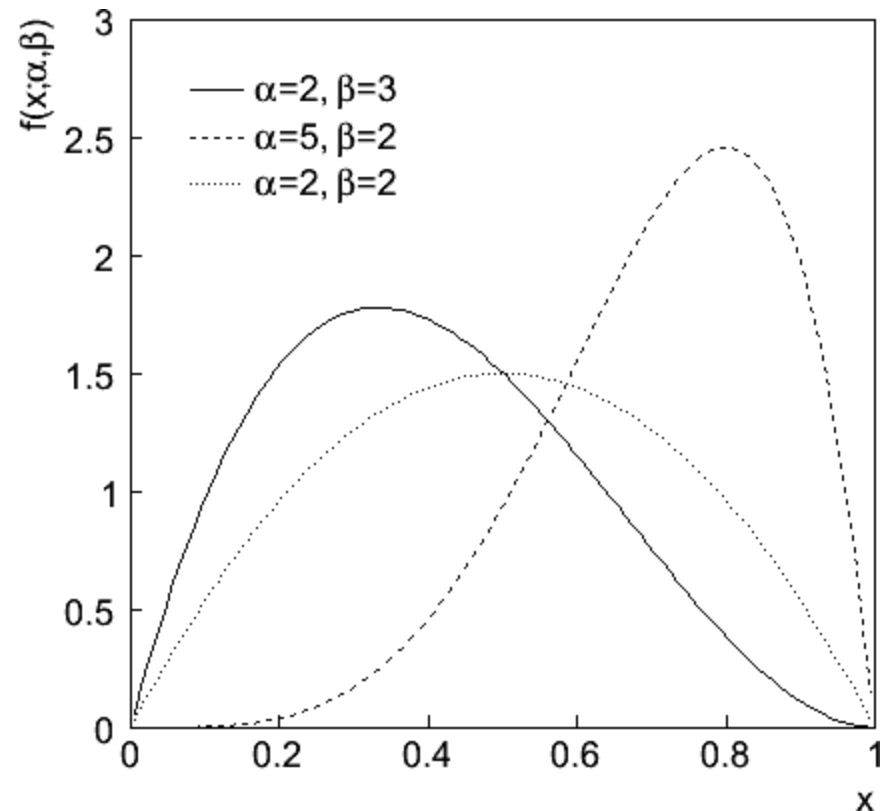
# Beta distribution

$$f(x; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1 - x)^{\beta-1}$$

$$E[x] = \frac{\alpha}{\alpha + \beta}$$

$$V[x] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

Often used to represent pdf of continuous r.v. nonzero only between finite limits.



# Gamma distribution

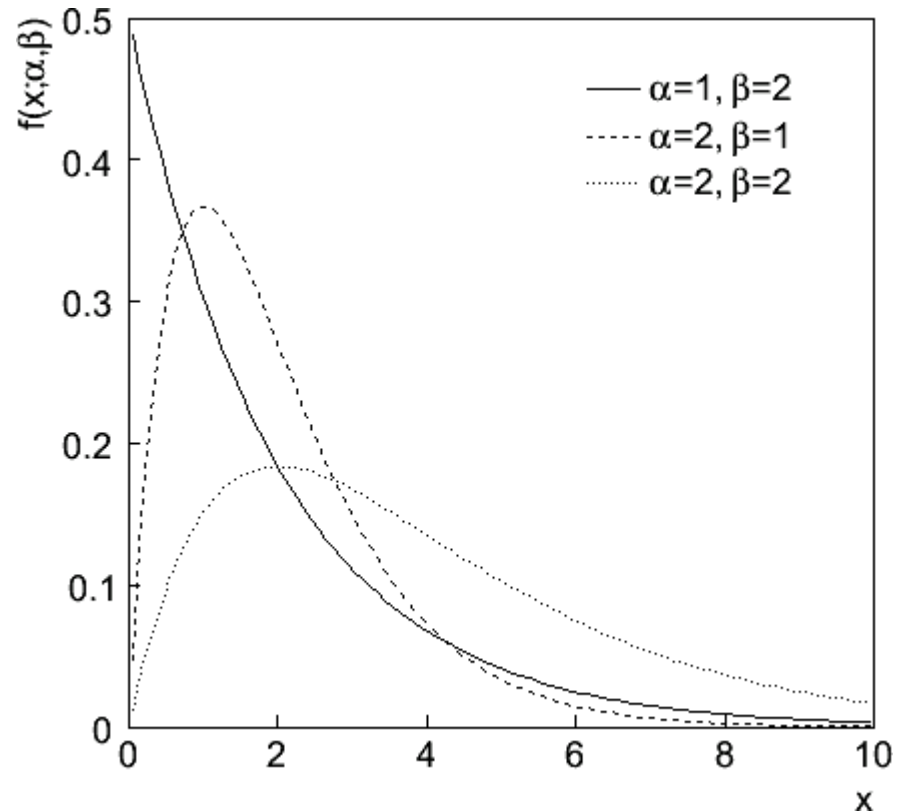
$$f(x; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}$$

$$E[x] = \alpha\beta$$

$$V[x] = \alpha\beta^2$$

Often used to represent pdf of continuous r.v. nonzero only in  $[0, \infty]$ .

Also e.g. sum of  $n$  exponential r.v.s or time until  $n$ th event in Poisson process  $\sim$  Gamma



# Student's $t$ distribution

$$f(x; \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi} \Gamma(\nu/2)} \left(1 + \frac{x^2}{\nu}\right)^{-\left(\frac{\nu+1}{2}\right)}$$

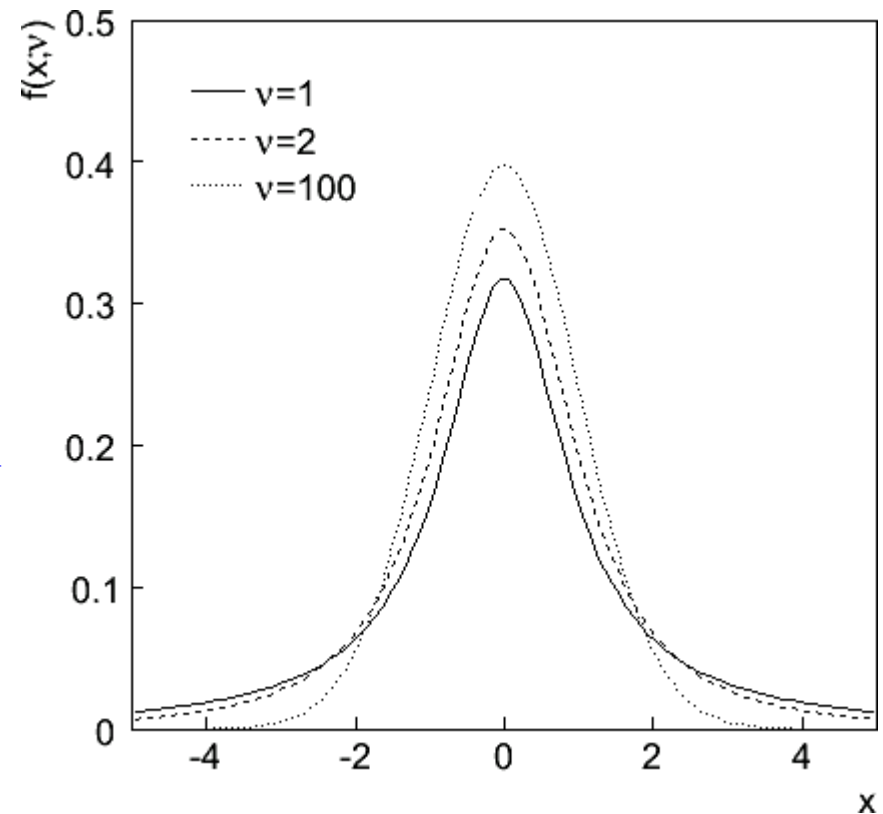
$$E[x] = 0 \quad (\nu > 1)$$

$$V[x] = \frac{\nu}{\nu - 2} \quad (\nu > 2)$$

$\nu$  = number of degrees of freedom  
(not necessarily integer)

$\nu = 1$  gives Cauchy,

$\nu \rightarrow \infty$  gives Gaussian.



## Student's $t$ distribution (2)

If  $x \sim$  Gaussian with  $\mu = 0$ ,  $\sigma^2 = 1$ , and

$z \sim \chi^2$  with  $n$  degrees of freedom, then

$t = x / (z/n)^{1/2}$  follows Student's  $t$  with  $\nu = n$ .

This arises in problems where one forms the ratio of a sample mean to the sample standard deviation of Gaussian r.v.s.

The Student's  $t$  provides a bell-shaped pdf with adjustable tails, ranging from those of a Gaussian, which fall off very quickly, ( $\nu \rightarrow \infty$ , but in fact already very Gauss-like for  $\nu =$  two dozen), to the very long-tailed Cauchy ( $\nu = 1$ ).

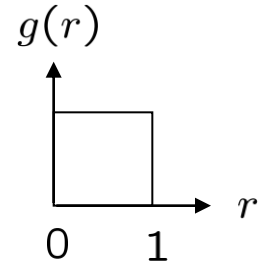
Developed in 1908 by William Gosset, who worked under the pseudonym "Student" for the Guinness Brewery.

# The Monte Carlo method

What it is: a numerical technique for calculating probabilities and related quantities using sequences of random numbers.

The usual steps:

- (1) Generate sequence  $r_1, r_2, \dots, r_m$  uniform in  $[0, 1]$ .
- (2) Use this to produce another sequence  $x_1, x_2, \dots, x_n$  distributed according to some pdf  $f(x)$  in which we're interested ( $x$  can be a vector).
- (3) Use the  $x$  values to estimate some property of  $f(x)$ , e.g., fraction of  $x$  values with  $a < x < b$  gives  $\int_a^b f(x) dx$ .
  - MC calculation = integration (at least formally)



MC generated values = ‘simulated data’

→ use for testing statistical procedures

# Random number generators

Goal: generate uniformly distributed values in  $[0, 1]$ .

Toss coin for e.g. 32 bit number... (too tiring).

→ ‘random number generator’

= computer algorithm to generate  $r_1, r_2, \dots, r_n$ .

Example: multiplicative linear congruential generator (MLCG)

$$n_{i+1} = (a n_i) \bmod m, \quad \text{where}$$

$$n_i = \text{integer}$$

$$a = \text{multiplier}$$

$$m = \text{modulus}$$

$$n_0 = \text{seed (initial value)}$$

N.B. mod = modulus (remainder), e.g.  $27 \bmod 5 = 2$ .

This rule produces a sequence of numbers  $n_0, n_1, \dots$

## Random number generators (2)

The sequence is (unfortunately) periodic!

Example (see Brandt Ch 4):  $a = 3, m = 7, n_0 = 1$

$$n_1 = (3 \cdot 1) \bmod 7 = 3$$

$$n_2 = (3 \cdot 3) \bmod 7 = 2$$

$$n_3 = (3 \cdot 2) \bmod 7 = 6$$

$$n_4 = (3 \cdot 6) \bmod 7 = 4$$

$$n_5 = (3 \cdot 4) \bmod 7 = 5$$

$$n_6 = (3 \cdot 5) \bmod 7 = 1 \quad \leftarrow \text{sequence repeats}$$

Choose  $a, m$  to obtain long period (maximum =  $m - 1$ );  $m$  usually close to the largest integer that can be represented in the computer.

Only use a subset of a single period of the sequence.

# Random number generators (3)

$r_i = n_i/m$  are in  $[0, 1]$  but are they ‘random’?

Choose  $a, m$  so that the  $r_i$  pass various tests of randomness:

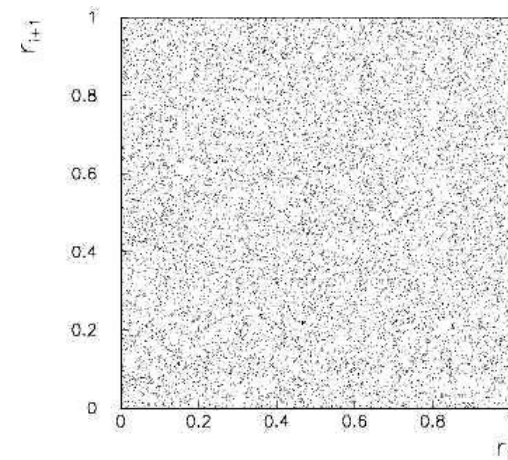
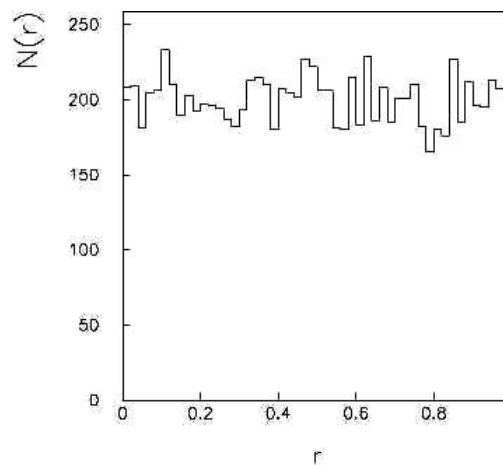
uniform distribution in  $[0, 1]$ ,

all values independent (no correlations between pairs),

e.g. L’Ecuyer, Commun. ACM **31** (1988) 742 suggests

$$a = 40692$$

$$m = 2147483399$$



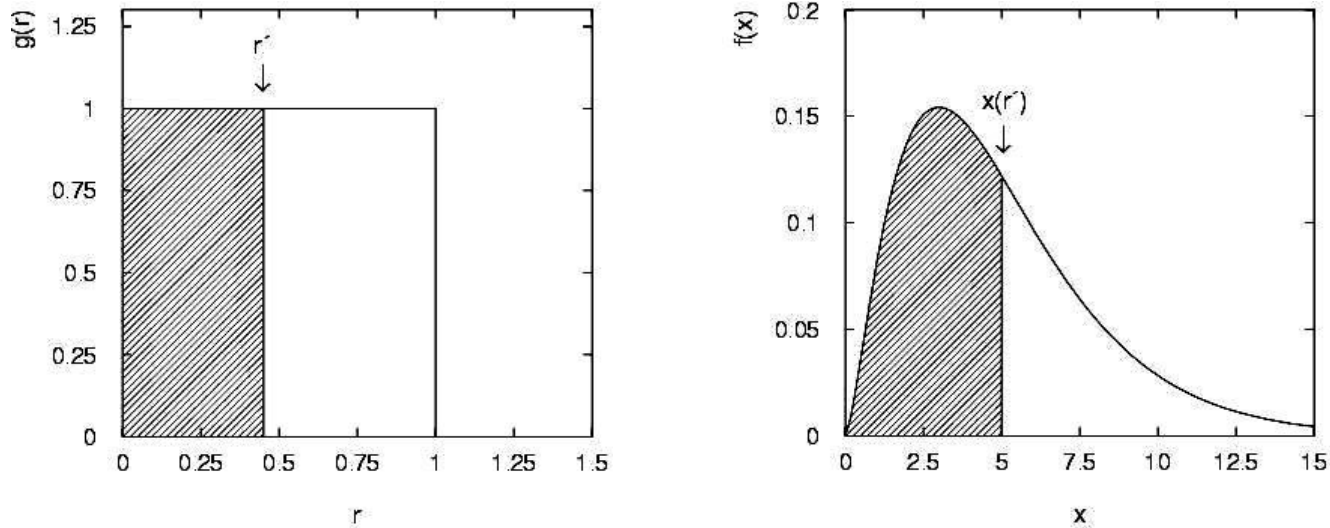
Far better algorithms available, e.g. **TRandom3**, period  $2^{19937} - 1$ .

See F. James, Comp. Phys. Comm. 60 (1990) 111; Brandt Ch. 4



# The transformation method

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



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

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

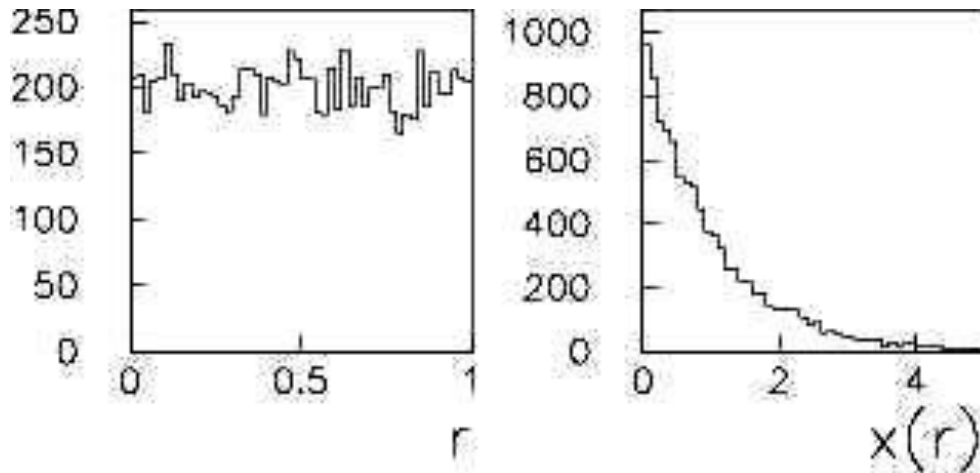
That is, set  $F(x) = r$  and solve for  $x(r)$ .

# Example of the transformation method

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

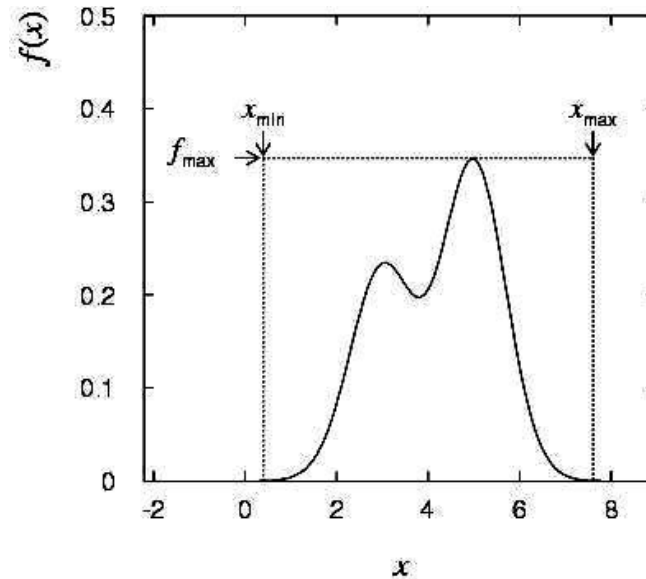
Set  $\int_0^x \frac{1}{\xi} e^{-x'/\xi} dx' = r$  and solve for  $x(r)$ .

→  $x(r) = -\xi \ln(1 - r)$  ( $x(r) = -\xi \ln r$  works too.)



# The acceptance-rejection method

Enclose the pdf in a box:



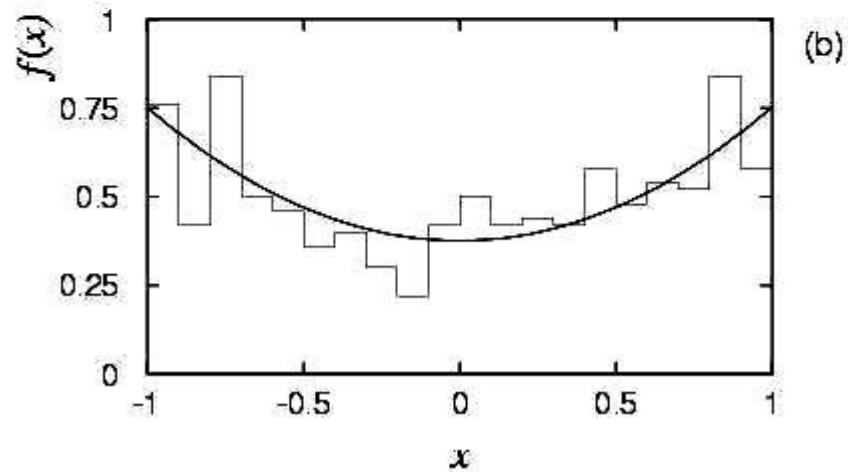
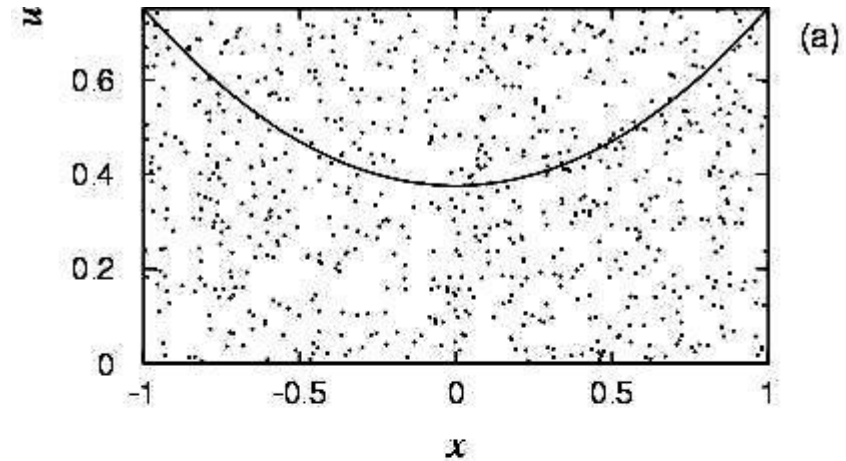
- (1) Generate a random number  $x$ , uniform in  $[x_{\min}, x_{\max}]$ , i.e.  
$$x = x_{\min} + r_1(x_{\max} - x_{\min})$$
,  $r_1$  is uniform in  $[0, 1]$ .
- (2) Generate a 2nd independent random number  $u$  uniformly distributed between 0 and  $f_{\max}$ , i.e.  $u = r_2 f_{\max}$ .
- (3) If  $u < f(x)$ , 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.



# Parameter estimation

The parameters of a pdf are constants that characterize its shape, e.g.

$$f(x; \theta) = \frac{1}{\theta} e^{-x/\theta}$$

r.v.                      parameter

Suppose we have a **sample** of observed values:  $\vec{x} = (x_1, \dots, x_n)$

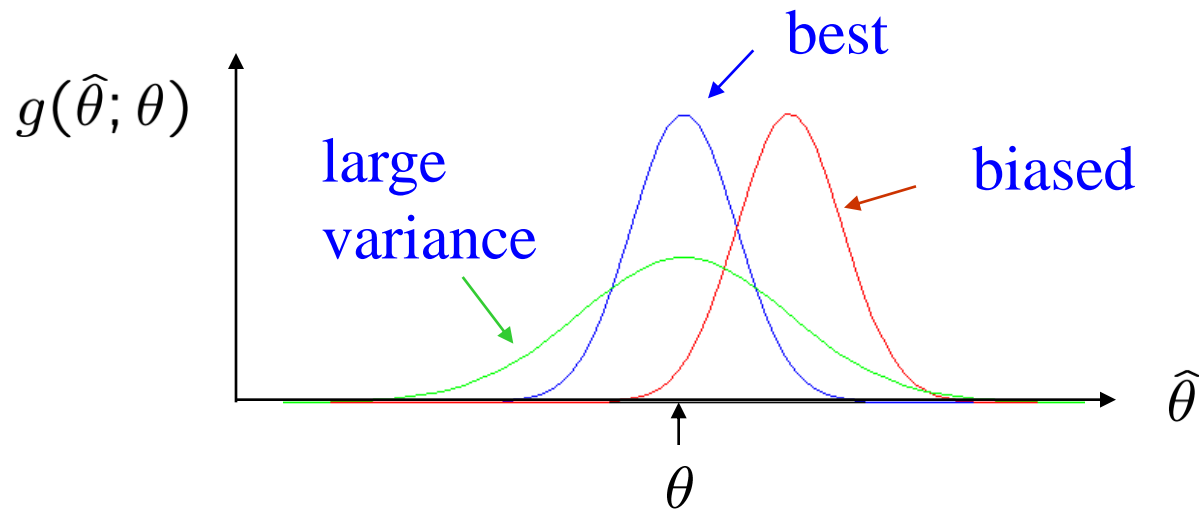
We want to find some function of the data to **estimate** the parameter(s):

$$\hat{\theta}(\vec{x}) \quad \leftarrow \text{estimator written with a hat}$$

Sometimes we say ‘estimator’ for the function of  $x_1, \dots, x_n$ ; ‘estimate’ for the value of the estimator with a particular data set.

# Properties of estimators

If we were to repeat the entire measurement, the estimates from each would follow a pdf:



We want small (or zero) bias (systematic error):  $b = E[\hat{\theta}] - \theta$

→ average of repeated measurements should tend to true value.

And we want a small variance (statistical error):  $V[\hat{\theta}]$

→ small bias & variance are in general conflicting criteria

# An estimator for the mean (expectation value)

Parameter:  $\mu = E[x]$

Estimator:  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \equiv \bar{x}$  ('sample mean')

We find:  $b = E[\hat{\mu}] - \mu = 0$

$$V[\hat{\mu}] = \frac{\sigma^2}{n} \quad \left( \sigma_{\hat{\mu}} = \frac{\sigma}{\sqrt{n}} \right)$$

# An estimator for the variance

Parameter:  $\sigma^2 = V[x]$

Estimator:  $\widehat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \equiv s^2$  ('sample variance')

We find:

$$b = E[\widehat{\sigma}^2] - \sigma^2 = 0 \quad (\text{factor of } n-1 \text{ makes this so})$$

$$V[\widehat{\sigma}^2] = \frac{1}{n} \left( \mu_4 - \frac{n-3}{n-1} \mu_2 \right), \quad \text{where}$$

$$\mu_k = \int (x - \mu)^k f(x) dx$$



# The likelihood function

Suppose the outcome of an experiment is:  $x_1, \dots, x_n$ , which is modeled as a sample from a joint pdf with parameter(s)  $\theta$ .

$$f(x_1, \dots, x_n; \theta)$$

Now evaluate this with the data sample obtained and regard it as a function of the parameter(s). This is the **likelihood function**:

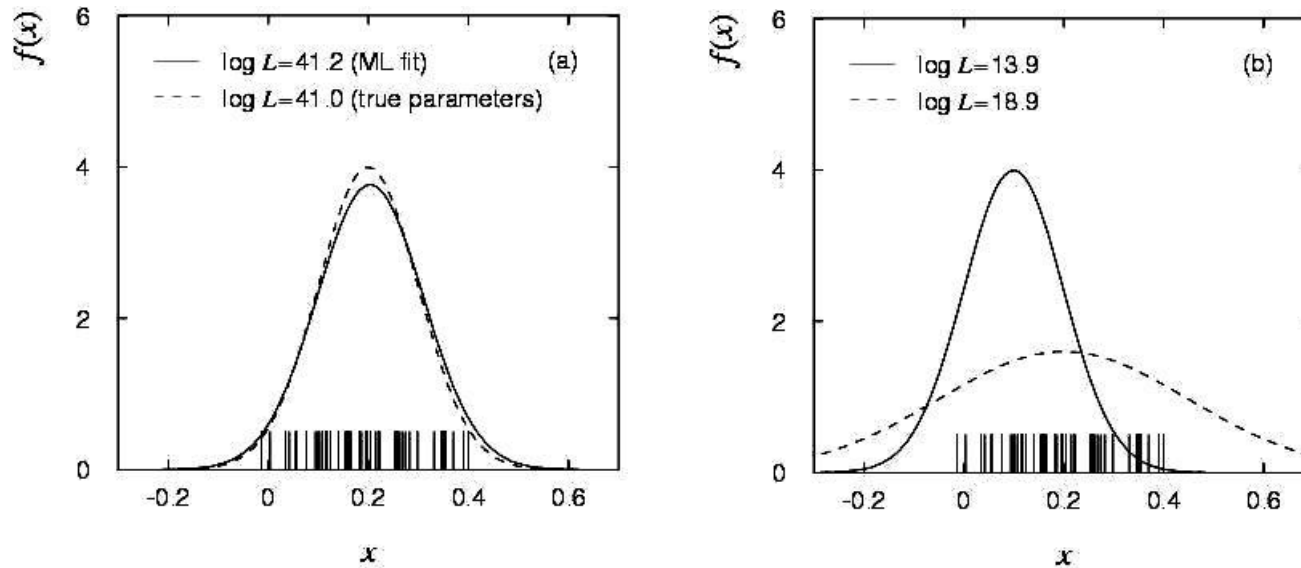
$$L(\theta) = f(x_1, \dots, x_n; \theta) \quad (x_i \text{ constant})$$

If the  $x_i$  are independent observations of  $x \sim f(x; \theta)$ , then,

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

# Maximum likelihood estimators

If the hypothesized  $\theta$  is close to the true value, then we expect a high probability to get data like that which we actually found.



So we define the maximum likelihood (ML) estimator(s) to be the parameter value(s) for which the likelihood is maximum.

ML estimators not guaranteed to have any ‘optimal’ properties, (but in practice they’re very good).

# ML example: parameter of exponential pdf

Consider exponential pdf,  $f(t; \tau) = \frac{1}{\tau} e^{-t/\tau}$

and suppose we have data,  $t_1, \dots, t_n$

The likelihood function is  $L(\tau) = \prod_{i=1}^n \frac{1}{\tau} e^{-t_i/\tau}$

The value of  $\tau$  for which  $L(\tau)$  is maximum also gives the maximum value of its logarithm (the log-likelihood function):

$$\ln L(\tau) = \sum_{i=1}^n \ln f(t_i; \tau) = \sum_{i=1}^n \left( \ln \frac{1}{\tau} - \frac{t_i}{\tau} \right)$$

# ML example: parameter of exponential pdf (2)

Find its maximum by setting  $\frac{\partial \ln L(\tau)}{\partial \tau} = 0$ ,

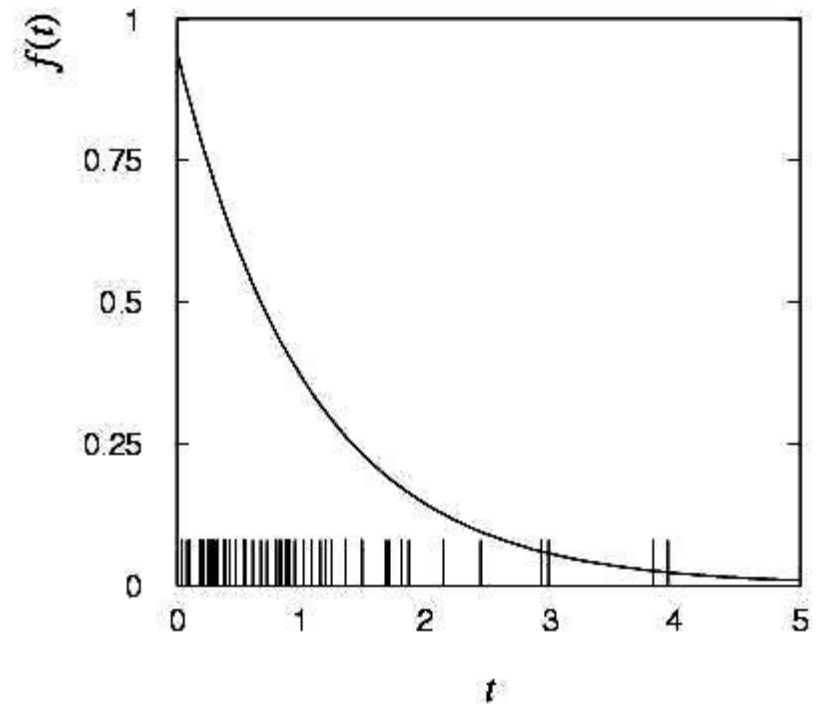
$$\rightarrow \hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$$

Monte Carlo test:

generate 50 values  
using  $\tau = 1$ :

We find the ML estimate:

$$\hat{\tau} = 1.062$$



# Variance of estimators: Monte Carlo method

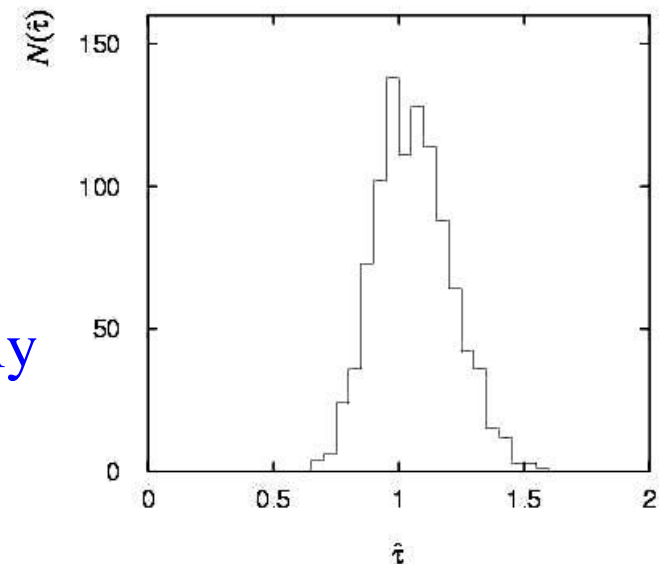
Having estimated our parameter we now need to report its ‘statistical error’, i.e., how widely distributed would estimates be if we were to repeat the entire measurement many times.

One way to do this would be to simulate the entire experiment many times with a Monte Carlo program (use ML estimate for MC).

For exponential example, from sample variance of estimates we find:

$$\hat{\sigma}_{\hat{\tau}} = 0.151$$

Note distribution of estimates is roughly Gaussian – (almost) always true for ML in large sample limit.



# Variance of estimators from information inequality

The **information inequality** (RCF) sets a lower bound on the variance of any estimator (not only ML):

$$V[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 \bigg/ E \left[ -\frac{\partial^2 \ln L}{\partial \theta^2} \right] \quad (b = E[\hat{\theta}] - \theta)$$

Often the bias  $b$  is small, and equality either holds exactly or is a good approximation (e.g. large data sample limit). Then,

$$V[\hat{\theta}] \approx -1 \bigg/ E \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right]$$

Estimate this using the 2nd derivative of  $\ln L$  at its maximum:

$$\hat{V}[\hat{\theta}] = - \left( \frac{\partial^2 \ln L}{\partial \theta^2} \right)^{-1} \bigg|_{\theta=\hat{\theta}}$$

# Variance of estimators: graphical method

Expand  $\ln L(\theta)$  about its maximum:

$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[ \frac{\partial \ln L}{\partial \theta} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2!} \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta})^2 + \dots$$

First term is  $\ln L_{\max}$ , second term is zero, for third term use information inequality (assume equality):

$$\ln L(\theta) \approx \ln L_{\max} - \frac{(\theta - \hat{\theta})^2}{2\widehat{\sigma}_{\hat{\theta}}^2}$$

$$\text{i.e., } \ln L(\hat{\theta} \pm \hat{\sigma}_{\hat{\theta}}) \approx \ln L_{\max} - \frac{1}{2}$$

→ to get  $\hat{\sigma}_{\hat{\theta}}$ , change  $\theta$  away from  $\hat{\theta}$  until  $\ln L$  decreases by 1/2.

# Example of variance by graphical method

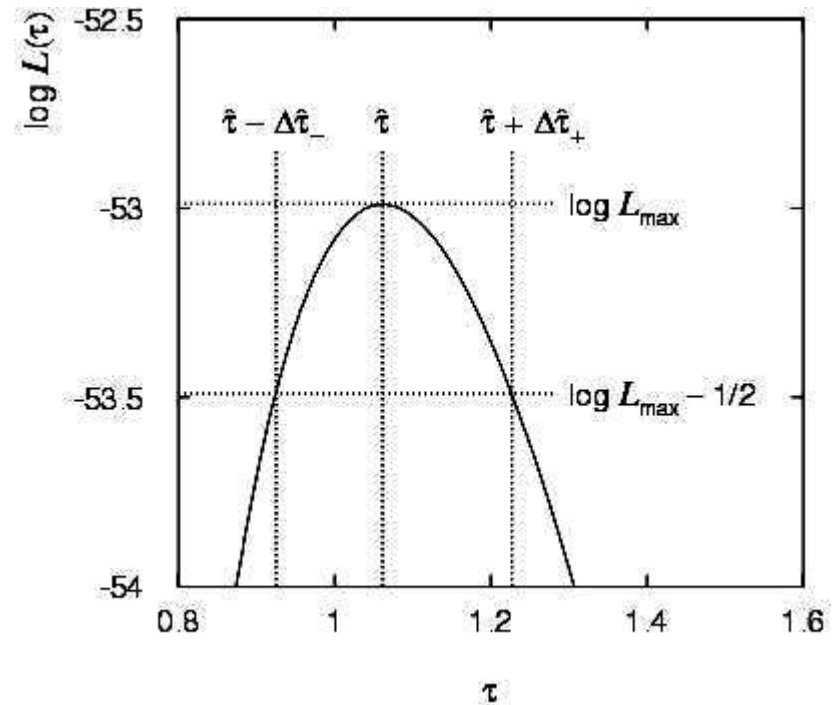
ML example with exponential:

$$\hat{\tau} = 1.062$$

$$\Delta\hat{\tau}_- = 0.137$$

$$\Delta\hat{\tau}_+ = 0.165$$

$$\hat{\sigma}_{\hat{\tau}} \approx \Delta\hat{\tau}_- \approx \Delta\hat{\tau}_+ \approx 0.15$$



Not quite parabolic  $\ln L$  since finite sample size ( $n = 50$ ).



# The method of least squares

Suppose we measure  $N$  values,  $y_1, \dots, y_N$ , assumed to be independent Gaussian r.v.s with

$$E[y_i] = \lambda(x_i; \theta) .$$

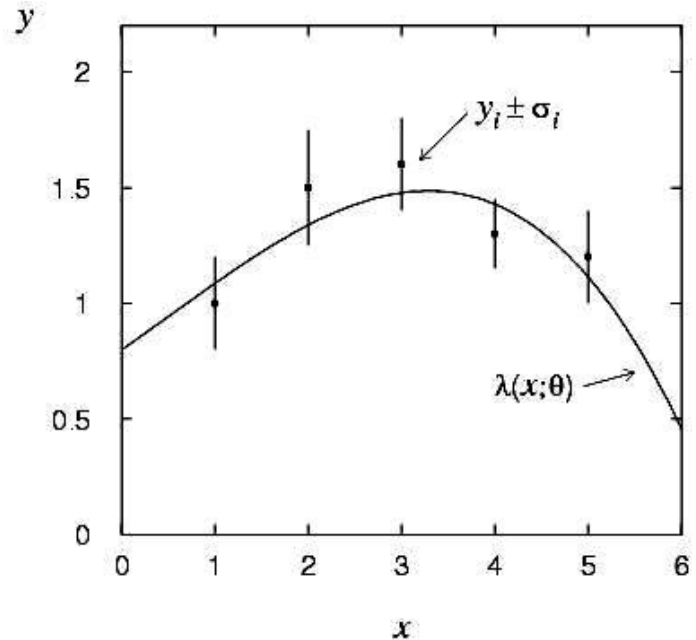
Assume known values of the control variable  $x_1, \dots, x_N$  and known variances

$$V[y_i] = \sigma_i^2 .$$

We want to estimate  $\theta$ , i.e., fit the curve to the data points.

The likelihood function is

$$L(\theta) = \prod_{i=1}^N f(y_i; \theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left[ -\frac{(y_i - \lambda(x_i; \theta))^2}{2\sigma_i^2} \right]$$



## The method of least squares (2)

The log-likelihood function is therefore

$$\ln L(\theta) = -\frac{1}{2} \sum_{i=1}^N \frac{(y_i - \lambda(x_i; \theta))^2}{\sigma_i^2} + \text{terms not depending on } \theta$$

So maximizing the likelihood is equivalent to minimizing

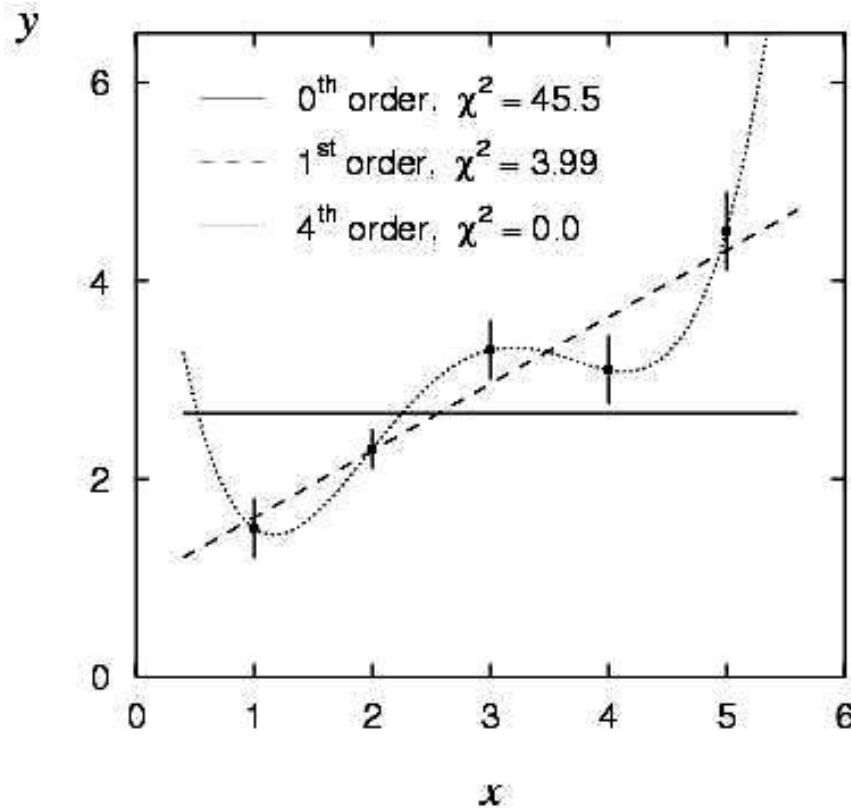
$$\chi^2(\theta) = \sum_{i=1}^N \frac{(y_i - \lambda(x_i; \theta))^2}{\sigma_i^2}$$

Minimum of this quantity defines the least squares estimator  $\hat{\theta}$ .

Often minimize  $\chi^2$  numerically (e.g. program **MINUIT**).

# Example of least squares fit

Fit a polynomial of order  $p$ :  $\lambda(x; \theta_0, \dots, \theta_p) = \sum_{n=0}^p \theta_n x^n$



# Variance of LS estimators

In most cases of interest we obtain the variance in a manner similar to ML. E.g. for data  $\sim$  Gaussian we have

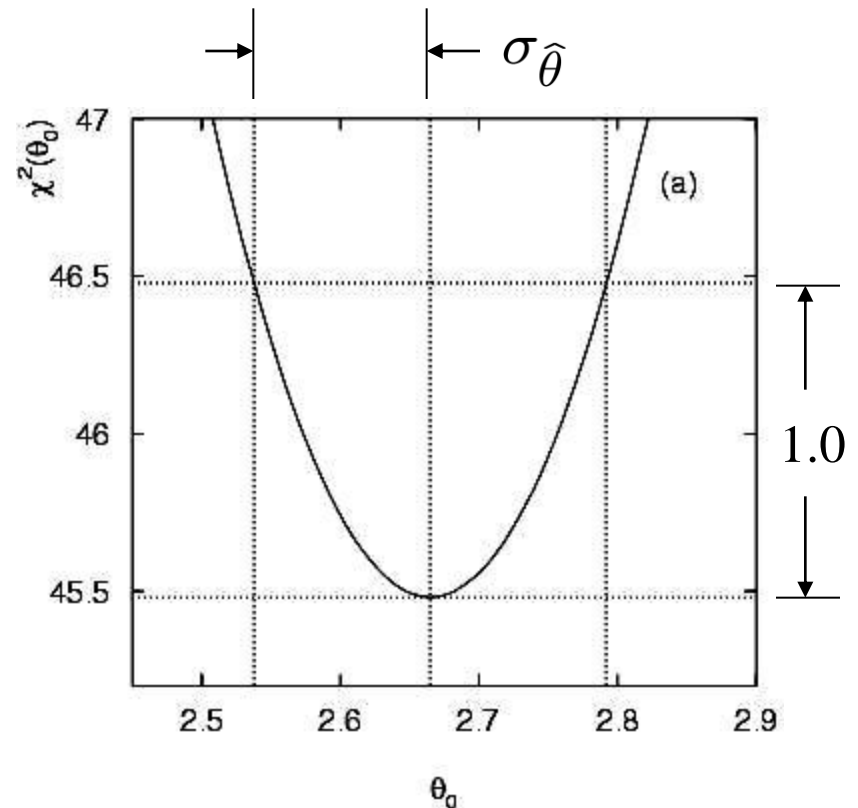
$$\chi^2(\theta) = -2 \ln L(\theta)$$

and so

$$\sigma_{\hat{\theta}}^2 \approx 2 \left[ \frac{\partial^2 \chi^2}{\partial \theta^2} \right]_{\theta=\hat{\theta}}^{-1}$$

or for the graphical method we take the values of  $\theta$  where

$$\chi^2(\theta) = \chi_{\min}^2 + 1$$



# Goodness-of-fit with least squares

The value of the  $\chi^2$  at its minimum is a measure of the level of agreement between the data and fitted curve:

$$\chi_{\min}^2 = \sum_{i=1}^N \frac{(y_i - \lambda(x_i; \hat{\theta}))^2}{\sigma_i^2}$$

It can therefore be employed as a goodness-of-fit statistic to test the hypothesized functional form  $\lambda(x; \theta)$ .

We can show that if the hypothesis is correct, then the statistic  $t = \chi_{\min}^2$  follows the chi-square pdf,

$$f(t; n_d) = \frac{1}{2^{n_d/2} \Gamma(n_d/2)} t^{n_d/2-1} e^{-t/2}$$

where the number of degrees of freedom is

$$n_d = \text{number of data points} - \text{number of fitted parameters}$$

## Goodness-of-fit with least squares (2)

The chi-square pdf has an expectation value equal to the number of degrees of freedom, so if  $\chi^2_{\min} \approx n_d$  the fit is ‘good’.

More generally, find the  $p$ -value: 
$$p = \int_{\chi^2_{\min}}^{\infty} f(t; n_d) dt$$

This is the probability of obtaining a  $\chi^2_{\min}$  as high as the one we got, or higher, if the hypothesis is correct.

E.g. for the previous example with 1st order polynomial (line),

$$\chi^2_{\min} = 3.99, \quad n_d = 5 - 2 = 3, \quad p = 0.263$$

whereas for the 0th order polynomial (horizontal line),

$$\chi^2_{\min} = 45.5, \quad n_d = 5 - 1 = 4, \quad p = 3.1 \times 10^{-9}$$

# Summary

We have quickly reviewed a large amount of material:

Probability

Distributions and their properties

Monte Carlo

Parameter estimation (ML, LS)

For a slower-paced treatment, see, e.g. the slides from the University of London course:

`www.pp.rhul.ac.uk/~cowan/stat\_course.html`

Next: statistical tests and multivariate methods