

# Introduction to Statistics – Day 2

## Lecture 1

Probability

Random variables, probability densities, etc.

## → Lecture 2

Brief catalogue of probability densities

The Monte Carlo method.

## Lecture 3

Statistical tests

Fisher discriminants, neural networks, etc

Significance and goodness-of-fit tests

## Lecture 4

Parameter estimation

Maximum likelihood and least squares

Interval estimation (setting limits)

# 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

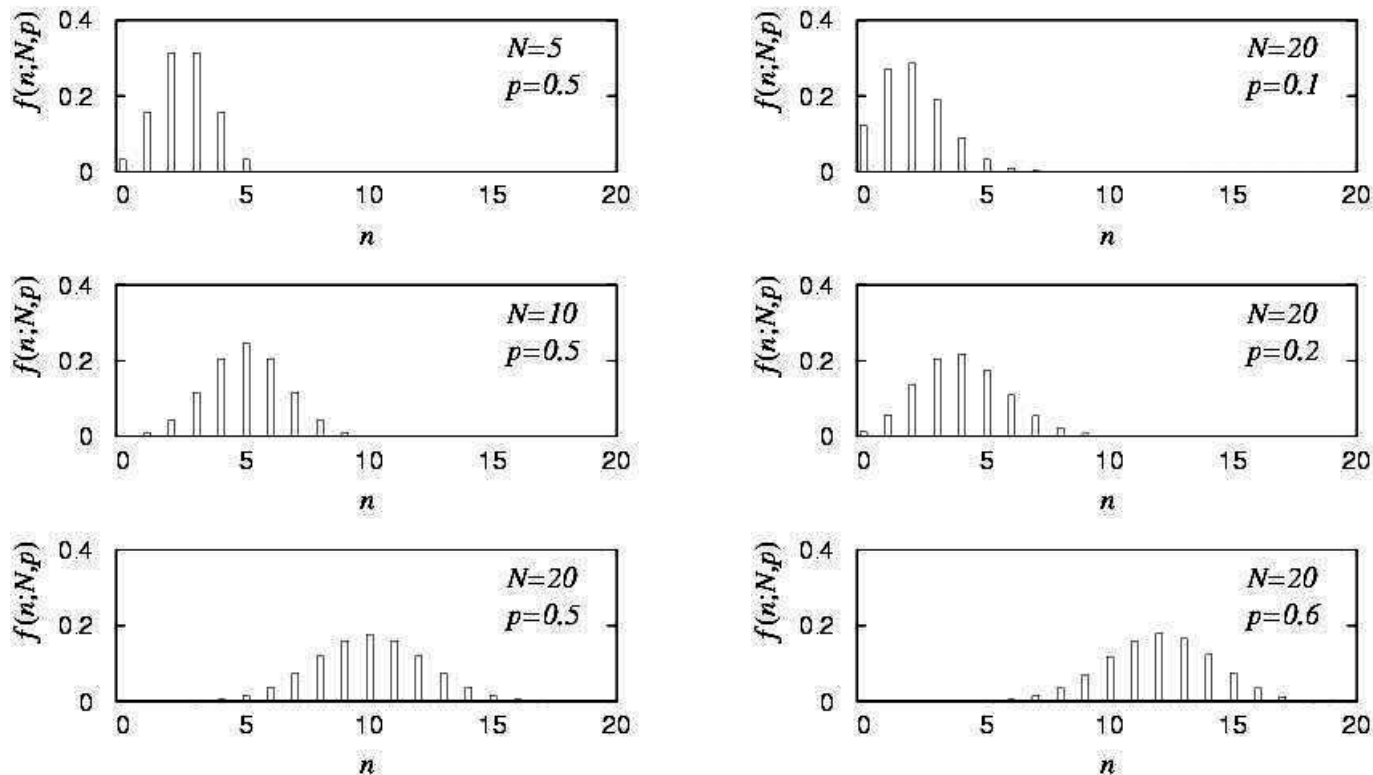
For the expectation value and variance we find:

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

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

# 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

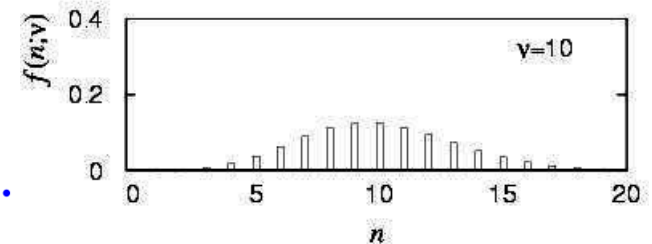
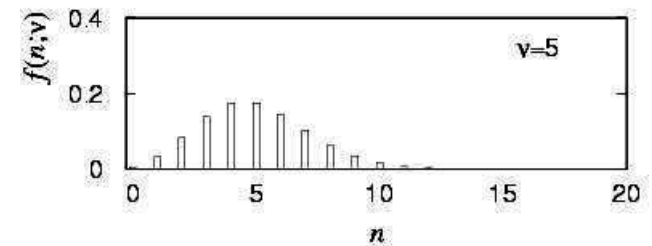
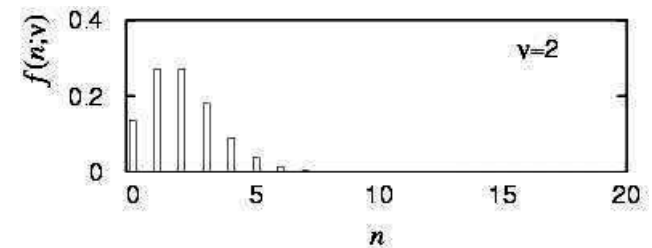
$$N \rightarrow \infty, \quad p \rightarrow 0, \quad E[n] = Np \rightarrow \nu .$$

→  $n$  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  $\nu = \sigma \int L dt$ .





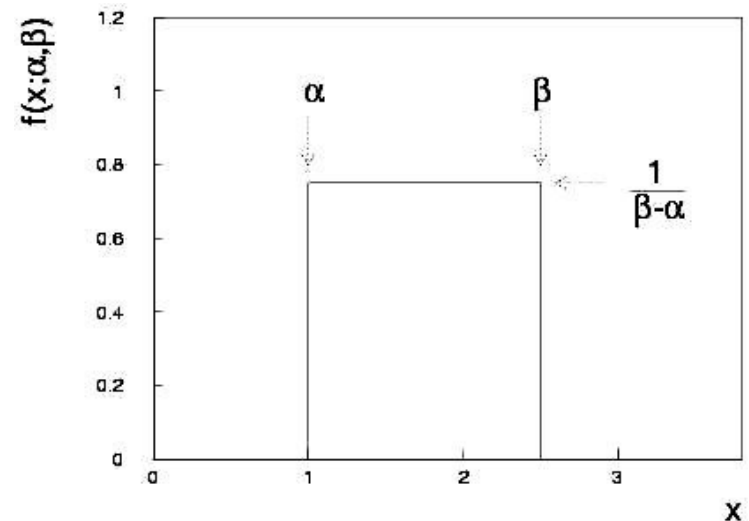
# Uniform distribution

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

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

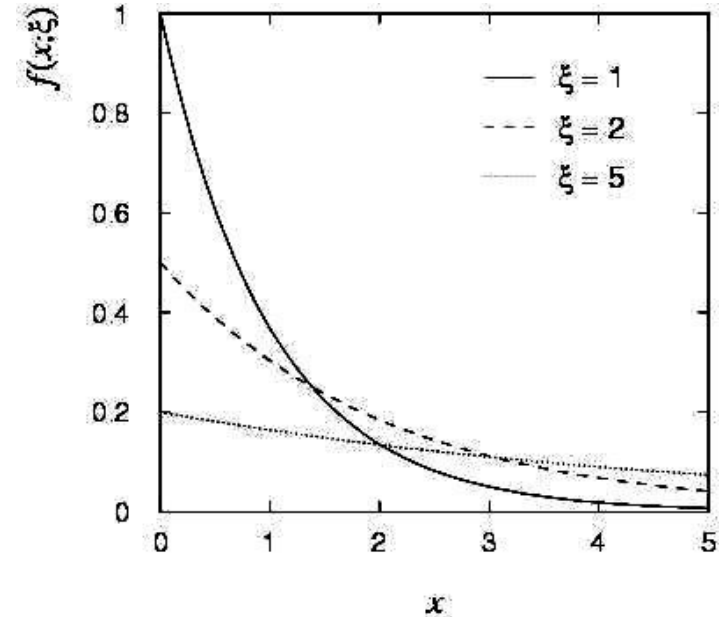
# Exponential distribution

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):  $f(t - t_0 | t \geq t_0) = f(t)$

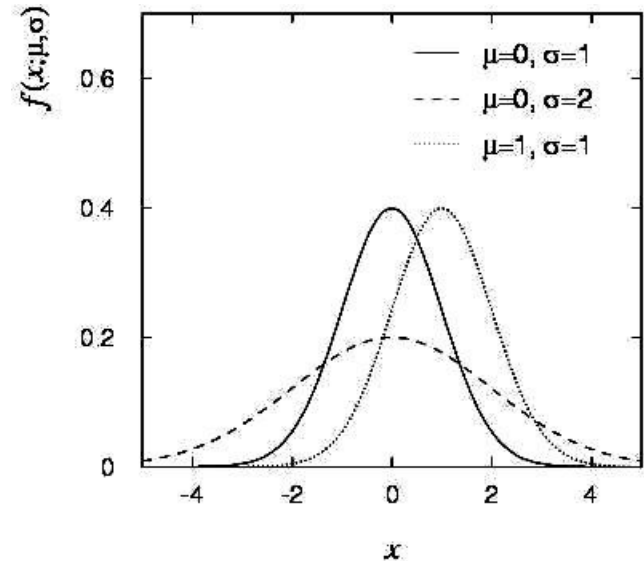
# Gaussian distribution

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  $y \sim \text{Gaussian}$  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.

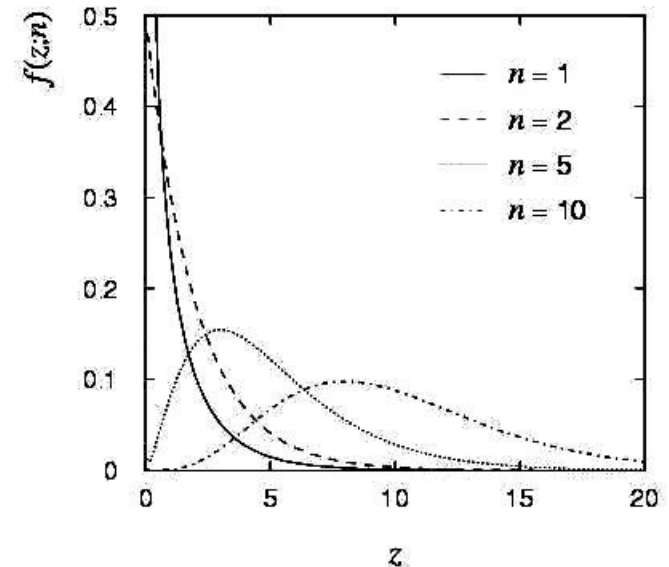
# Chi-square ( $\chi^2$ ) distribution

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

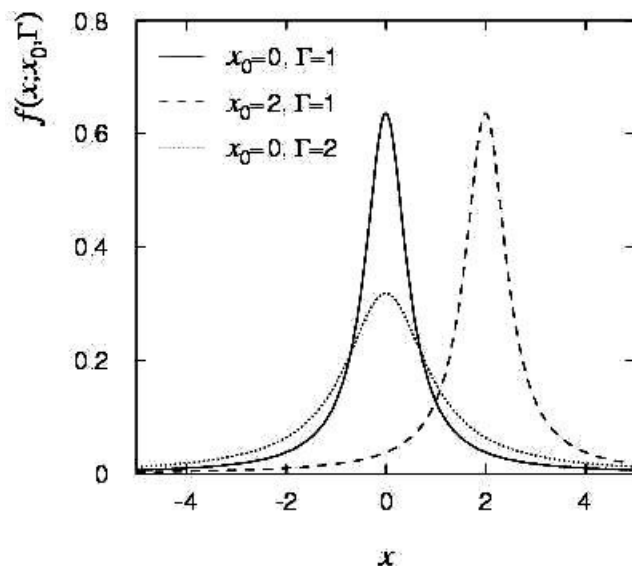
$$f(x; \Gamma, x_0) = \frac{1}{\pi} \frac{\Gamma/2}{\Gamma^2/4 + (x - x_0)^2}$$

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

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

$x_0$  = mode (most probable value)

$\Gamma$  = full width at half maximum



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

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



# Landau distribution

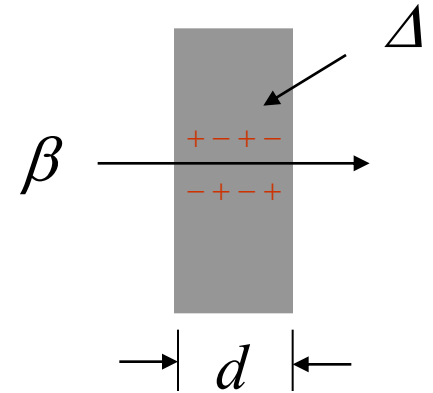
For a charged particle with  $\beta = v/c$  traversing a layer of matter of thickness  $d$ , the energy loss  $\Delta$  follows the Landau pdf:

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

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

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

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



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

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

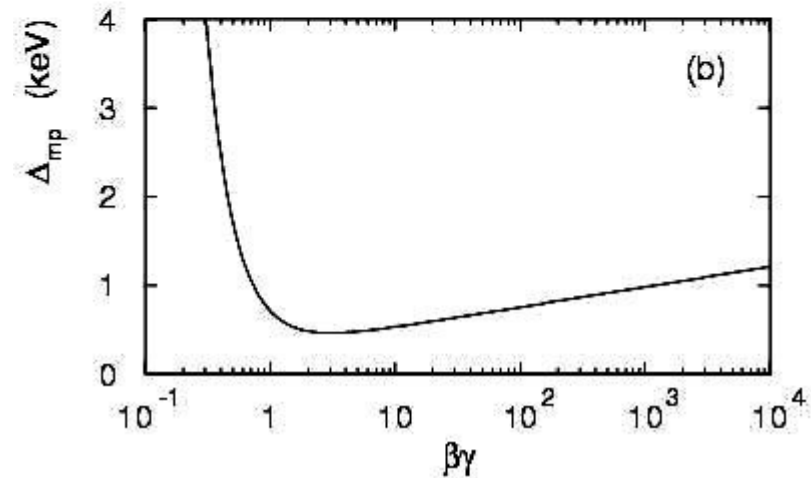
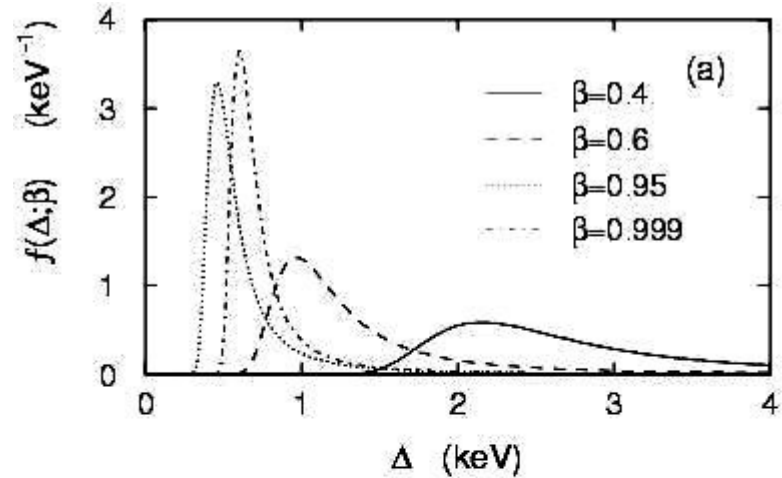
# Landau distribution (2)

Long ‘Landau tail’

→ all moments  $\infty$

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

→ particle i.d.



# The Monte Carlo method

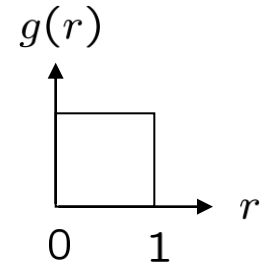
What it is: a numerical technique for calculating probabilities and related quantities using sequences of random numbers.

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)

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

$n_i$  = integer

$a$  = multiplier

$m$  = modulus

$n_0$  = seed (initial value)

N.B.  $\bmod$  = 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:

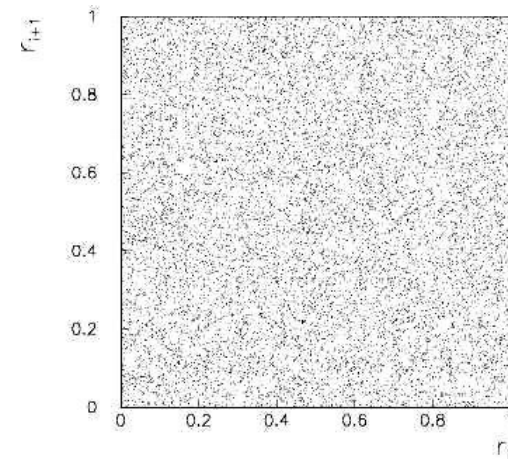
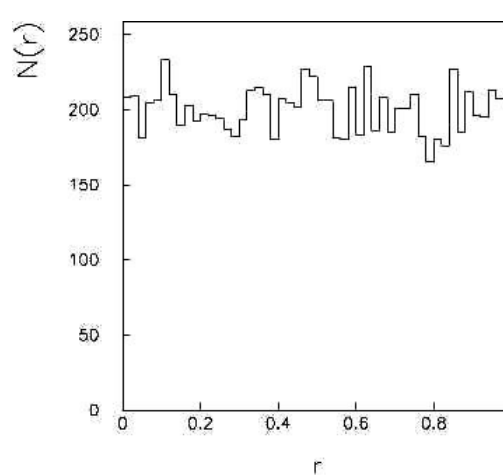
uniform distribution in  $[0, 1]$ ,

all values independent (no correlations between pairs),

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

$$a = 40692$$

$$m = 2147483399$$

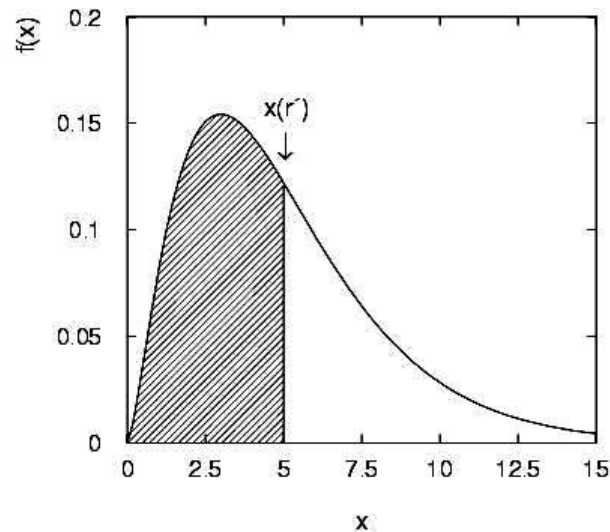
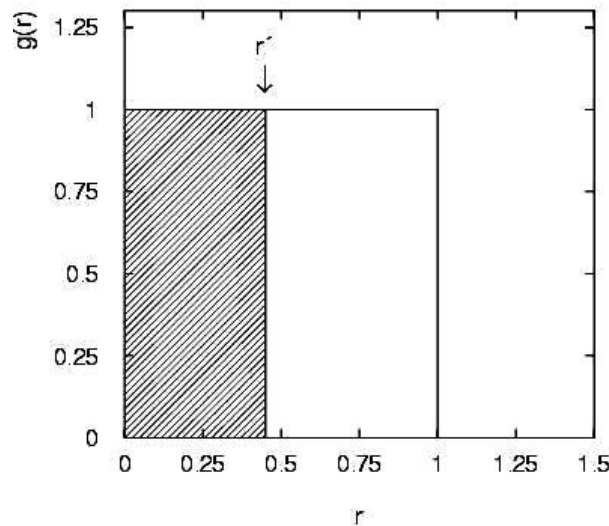


Far better algorithms available, e.g. **TRandom3**, period  $2^{19937} - 1$ .

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

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

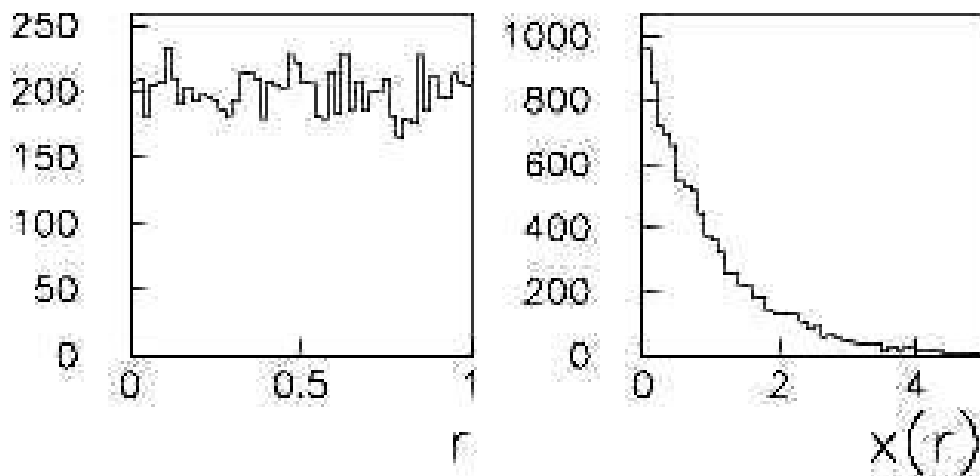
That is, set  $F(x) = r$  and solve for  $x(r)$ .

# 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  $x(r)$ .

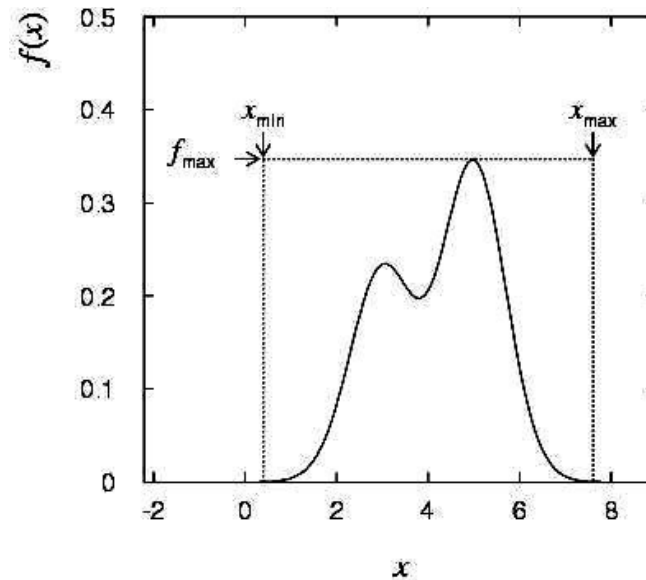
→  $x(r) = -\xi \ln(1 - r)$  ( $x(r) = -\xi \ln r$  works too.)





# The acceptance-rejection method

Enclose the pdf in a box:



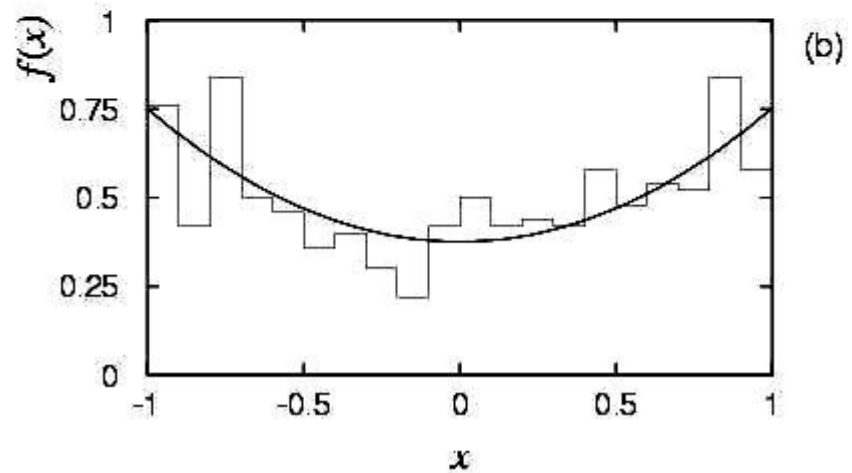
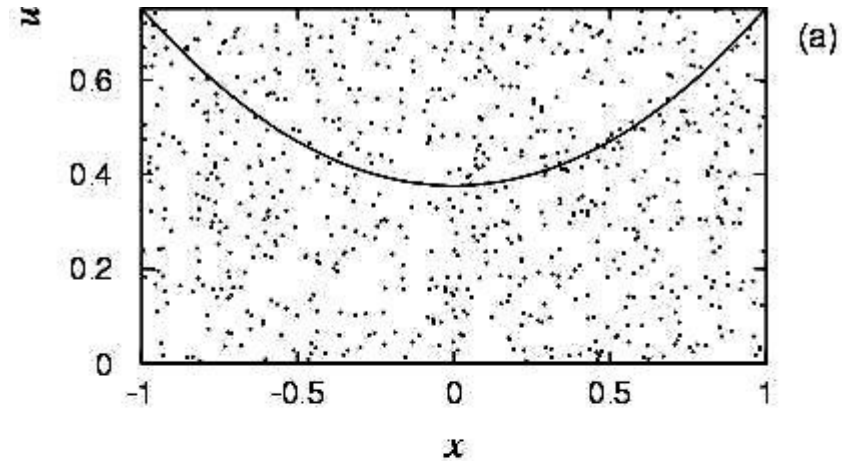
- (1) Generate a random number  $x$ , uniform in  $[x_{\min}, x_{\max}]$ , i.e.  
$$x = x_{\min} + r_1(x_{\max} - x_{\min})$$
,  $r_1$  is uniform in  $[0,1]$ .
- (2) Generate a 2nd independent random number  $u$  uniformly distributed between 0 and  $f_{\max}$ , i.e.  $u = r_2 f_{\max}$ .
- (3) If  $u < f(x)$ , then accept  $x$ . If not, reject  $x$  and repeat.

# Example with acceptance-rejection method

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

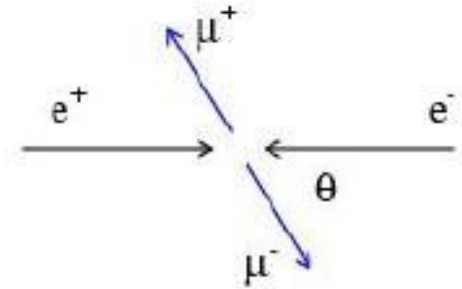
$$(-1 \leq x \leq 1)$$

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



# Monte Carlo event generators

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



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

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

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

Less simple: ‘event generators’ for a variety of reactions:

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

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

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

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

# A simulated event

I	particle/jet	KS	KF	orig	p_x	p_y	p_z	E	m
1	!p+	21	2212	0	0.000	0.000	7000.000	7000.000	0.938
2	!p+	21	2212	0	0.000	0.000	-7000.000	7000.000	0.938
=====									
3	!g!	21	21	1	0.863	-0.323	1739.862	1739.862	0.000
4	!ubar!	21	-2	2	-0.621	-0.163	-777.415	777.415	0.000
5	!g!	21	21	3	-2.427	5.486	1487.857	1487.857	0.000
6	!g!	21	21	4	-62.910	63.357	-463.274	471.274	0.000
7	!~g!	21	1000021	0	314.363	544.843	498.897	979.897	0.000
8	!~g!	21	1000021	0	-379.700	-476.000	525.686	980.686	0.000
9	!~chi_1-!	21	-1000024	7	130.058	112.247	129.860	263.860	0.000
10	!sbar!	21	-3	7	259.400	187.468	83.100	330.100	0.000
11	!c!	21	4	7	-79.403	242.409	283.026	381.026	0.000
12	!~chi_20!	21	1000023	8	-326.241	-80.971	113.712	385.712	0.000
13	!b!	21	5	8	-51.841	-294.077	389.853	491.853	0.000
14	!bbar!	21	-5	8	-0.597	-99.577	21.299	101.299	0.000
15	!~chi_10!	21	1000022	9	103.352	81.316	83.457	175.457	0.000
16	!s!	21	3	9	5.451	38.374	52.302	65.302	0.000
17	!cbar!	21	-4	9	20.839	-7.250	-5.938	22.938	0.000
18	!~chi_10!	21	1000022	12	-136.266	-72.961	53.246	181.246	0.000
19	!nu_mu!	21	14	12	-78.263	-24.757	21.719	84.719	0.000
20	!nu_mubar!	21	-14	12	-107.801	16.901	38.226	115.226	0.000
=====									
21	gamma	1	22	4	2.636	1.357	0.125	2.636	0.000
22	(~chi_1-)	11	-1000024	9	129.643	112.440	129.820	262.820	0.000
23	(~chi_20)	11	1000023	12	-322.330	-80.817	113.191	382.191	0.000
24	(~chi_10)	1	1000022	15	97.944	77.819	80.917	169.917	0.000
25	(~chi_10)	1	1000022	18	-136.266	-72.961	53.246	181.246	0.000
26	nu_mu	1	14	19	-78.263	-24.757	21.719	84.719	0.000
27	nu_mubar	1	-14	20	-107.801	16.901	38.226	115.226	0.000
28	(Delta++)	11	2224	2	0.222	0.012	-2734.287	2734.287	0.000

397	pi+	1	211	209	0.006	0.398	-308.296	308.297	0.140
398	gamma	1	22	211	0.407	0.087	-1695.458	1695.458	0.000
399	gamma	1	22	211	0.113	-0.029	-314.822	314.822	0.000
400	(pi0)	11	111	212	0.021	0.122	-103.709	103.709	0.135
401	(pi0)	11	111	212	0.084	-0.068	-94.276	94.276	0.135
402	(pi0)	11	111	212	0.267	-0.052	-144.673	144.674	0.135
403	gamma	1	22	215	-1.581	2.473	3.306	4.421	0.000
404	gamma	1	22	215	-1.494	2.143	3.051	4.016	0.000
405	pi-	1	-211	216	0.007	0.738	4.015	4.085	0.140
406	pi+	1	211	216	-0.024	0.293	0.486	0.585	0.140
407	K+	1	321	218	4.382	-1.412	-1.799	4.968	0.494
408	pi-	1	-211	218	1.183	-0.894	-0.176	1.500	0.140
409	(pi0)	11	111	218	0.955	-0.459	-0.590	1.221	0.135
410	(pi0)	11	111	218	2.349	-1.105	-1.181	2.855	0.135
411	(Kbar0)	11	-311	219	1.441	-0.247	-0.472	1.615	0.498
412	pi-	1	-211	219	2.232	-0.400	-0.249	2.285	0.140
413	K+	1	321	220	1.380	-0.652	-0.361	1.644	0.494
414	(pi0)	11	111	220	1.078	-0.265	0.175	1.132	0.135
415	(K_S0)	11	310	222	1.841	0.111	0.894	2.109	0.498
416	K+	1	321	223	0.307	0.107	0.252	0.642	0.494
417	pi-	1	-211	223	0.266	0.316	-0.201	0.480	0.140
418	nbar0	1	-2112	226	1.335	1.641	2.078	3.111	0.940
419	(pi0)	11	111	226	0.899	1.046	1.311	1.908	0.135
420	pi+	1	211	227	0.217	1.407	1.356	1.971	0.140
421	(pi0)	11	111	227	1.207	2.336	2.767	3.820	0.135
422	n0	1	2112	228	3.475	5.324	5.702	8.592	0.940
423	pi-	1	-211	228	1.856	2.606	2.808	4.259	0.140
424	gamma	1	22	229	-0.012	0.247	0.421	0.489	0.000
425	gamma	1	22	229	0.025	0.034	0.009	0.043	0.000
426	pi+	1	211	230	2.718	5.229	6.403	8.703	0.140
427	(pi0)	11	111	230	4.109	6.747	7.597	10.961	0.135
428	pi-	1	-211	231	0.551	1.233	1.945	2.372	0.140
429	(pi0)	11	111	231	0.645	1.141	0.922	1.608	0.135
430	gamma	1	22	232	-0.383	1.169	1.208	1.724	0.000
431	gamma	1	22	232	-0.201	0.070	0.060	0.221	0.000

PYTHIA Monte Carlo  
pp → gluino-gluino

# Monte Carlo detector simulation

Takes as input the particle list and momenta from generator.

Simulates detector response:

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

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

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

Estimate ‘efficiencies’ = #events found / # events generated.

Programming package: **GEANT**

## Wrapping up lecture 2

We've looked at a number of important distributions:

Binomial, Multinomial, Poisson, Uniform, Exponential  
Gaussian, Chi-square, Cauchy, Landau,

and we've seen the Monte Carlo method:

calculations based on sequences of random numbers,  
used to simulate particle collisions, detector response.

So far, we've mainly been talking about **probability**.

But suppose now we are faced with experimental data.

We want to infer something about the (probabilistic) processes  
that produced the data.

This is **statistics**, the main subject of the next two lectures.