# Statistical Methods
# in High Energy Physics -1

# Outline

- Probability
- Random variables, probability densities

- Probability densities in HEP

- Parameter estimation, hypothesis tests
- Maximum likelihood and least squares

Ref: Detection and Estimation Theory by Van Trees;
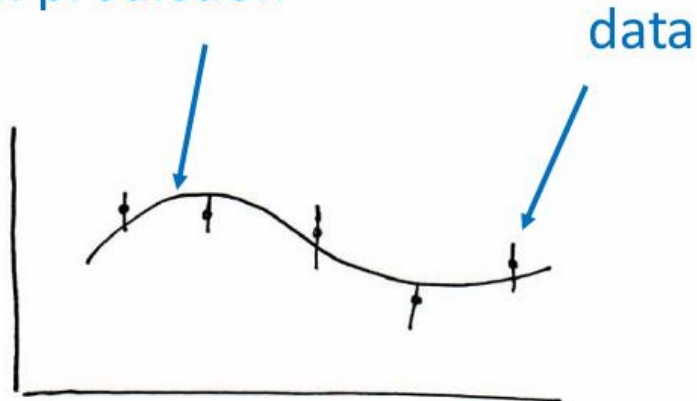    Statistical data analysis by Cowan

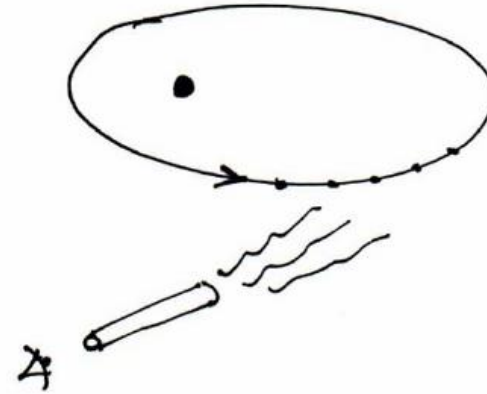# Theory ↔ Statistics ↔ Experiment

Theory (model, hypothesis):

Experiment (observation):

$$F = - G \frac{m_1 m_2}{r^2} \quad , \quad \cdot \cdot \cdot$$

+ response of measurement apparatus

= model prediction
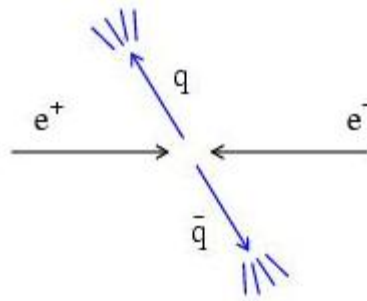
data



Uncertainty enters on many levels

→ quantify with probability

# 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?)
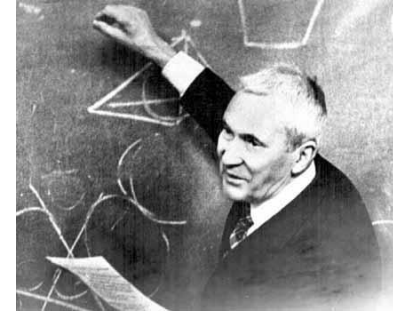
# Definition of probability



Consider a set *S* with subsets *A*, *B*, ...

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

$$P(S) = 1$$

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

Kolmogorov axioms (1933)

From these axioms we can derive further properties, e.g.

$$P(\overline{A}) = 1 - P(A)$$

$$P(A \cup \overline{A}) = 1$$

$$P(\emptyset) = 0$$

$$\text{if } A \subset B, \text{ then } P(A) \leq P(B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

# Conditional probability, independence

Also define conditional probability of *A* given *B* (with *P(B) ≠ 0*):

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

E.g. rolling dice: $P(n < 3 \mid n \text{ even}) = \frac{P((n<3) \cap n \text{ even})}{P(\text{even})} = \frac{1/6}{3/6} = \frac{1}{3}$

Subsets *A, B* independent if: $P(A \cap B) = P(A)P(B)$

If *A, B* independent, $P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A)$

do not confuse with disjoint subsets, i.e., $A \cap B = \emptyset$

# Interpretation of probability

I. **Relative frequency**

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

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

quantum mechanics, particle scattering, radioactive decay...

II. **Subjective probability (Bayesian)**

$A, B, \ldots$ 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.

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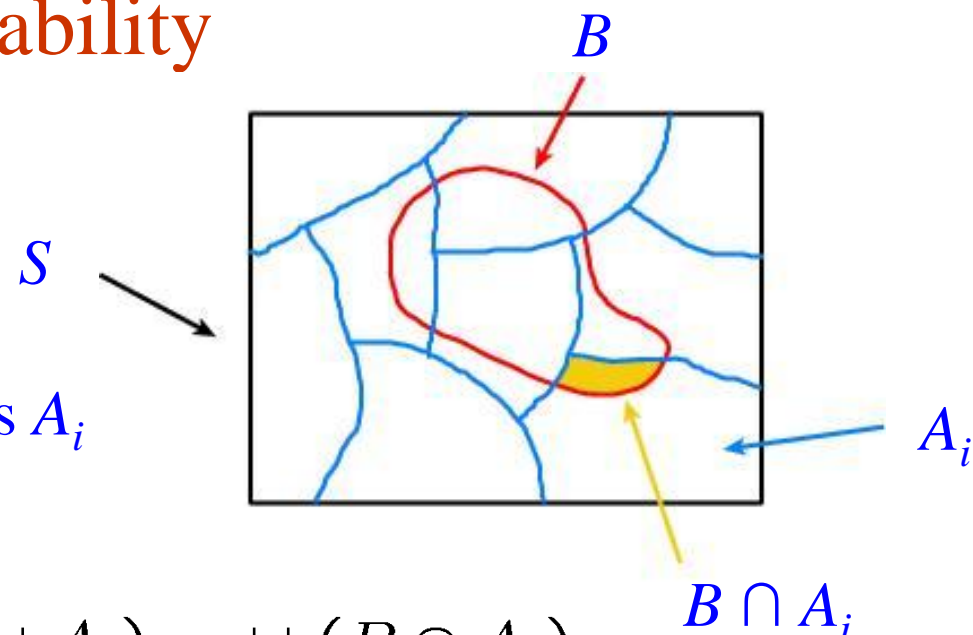
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

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*

*B*

$A_i$

$B \cap A_i$

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

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

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

Bayes' theorem becomes

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

# Bayesian Statistics – general philosophy

In Bayesian statistics, use subjective probability for hypotheses:

probability of the data assuming
hypothesis $H$ (the likelihood)

prior probability, i.e.,
before seeing the data

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

posterior probability, i.e.,
after seeing the data

normalization involves sum
over all possible hypotheses

Bayes' theorem has an "if-then" character:  If your prior
probabilities were $\pi(H)$, then it says how these probabilities
should change in the light of the data.

No unique prescription for priors (subjective!)

# An example using Bayes' theorem

Suppose the probability (for anyone) to have a disease D is:

$$P(\text{D}) = 0.001$$

$$P(\text{no D}) = 0.999$$

$\leftarrow$ prior probabilities, i.e., before any test carried out

Consider a test for the disease: result is + or −

$$P(+|\text{D}) = 0.98$$

$$P(-|\text{D}) = 0.02$$

$\leftarrow$ probabilities to (in)correctly identify a person with the disease

$$P(+|\text{no D}) = 0.03$$

$$P(-|\text{no D}) = 0.97$$

$\leftarrow$ probabilities to (in)correctly identify a healthy person

Suppose your result is +. How worried should you be?

# Bayes' theorem example (cont.)

The probability to have the disease given a + result is

$$p(\text{D}|+) = \frac{P(+|\text{D})P(\text{D})}{P(+|\text{D})P(\text{D}) + P(+|\text{no D})P(\text{no D})}$$

$$= \frac{0.98 \times 0.001}{0.98 \times 0.001 + 0.03 \times 0.999}$$

$$= 0.032 \qquad \leftarrow \text{posterior probability}$$

i.e. you're probably OK!

Your viewpoint:  my degree of belief that I have the disease is 3.2%.

Your doctor's viewpoint:  3.2% of people like this have the disease.

# 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)

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

Or for discrete outcome $x_i$ with e.g. $i = 1, 2, ...$ we have
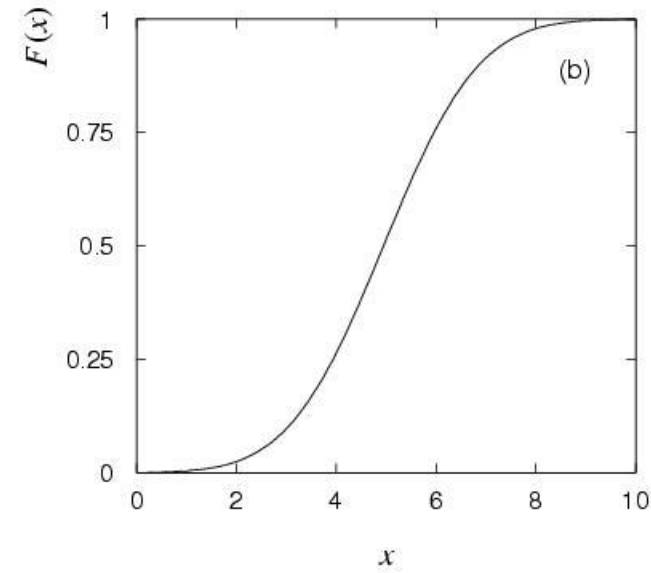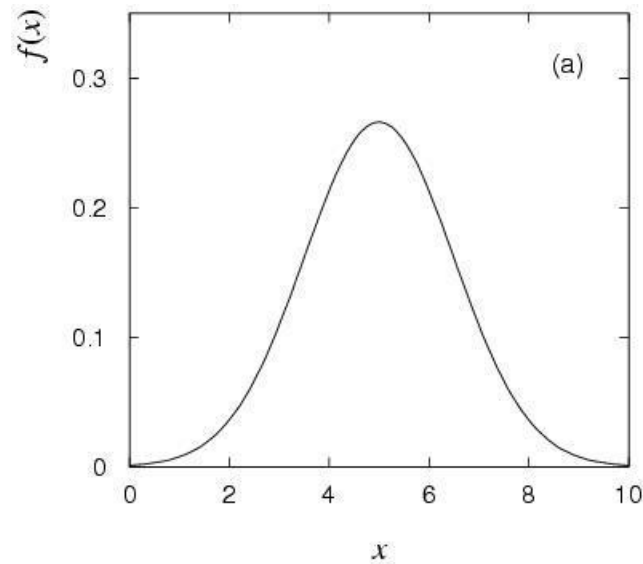
$$P(x_i) = p_i \qquad \text{probability mass function}$$

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

# Cumulative distribution function

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

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



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

# Other types of probability densities

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

$\rightarrow$ joint pdf  $f(x_1, \ldots, x_n)$

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

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

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

Sometimes we want to consider some components as constant

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

# Expectation values

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

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

Notation (often): $\quad E[x] = \mu \quad$ ~ "centre of gravity" of pdf.
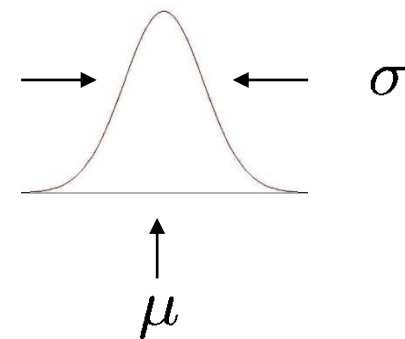
For a function $y(x)$ with pdf $g(y)$,

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

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

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

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

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

# Covariance and correlation

Define covariance 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

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

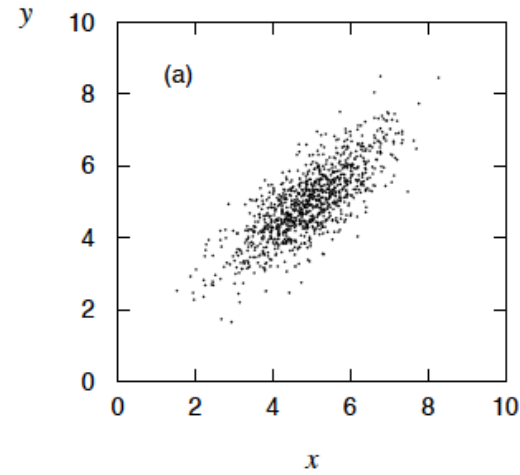If *x*, *y*, independent, i.e., $f(x, y) = f_x(x) f_y(y)$, then

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

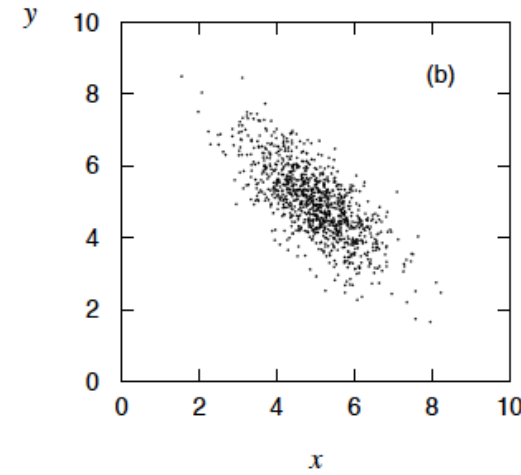$\rightarrow \quad \text{cov}[x, y] = 0 \qquad$ *x* and *y*, 'uncorrelated'
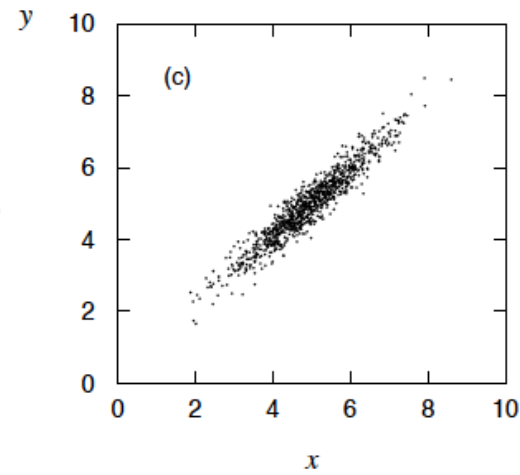
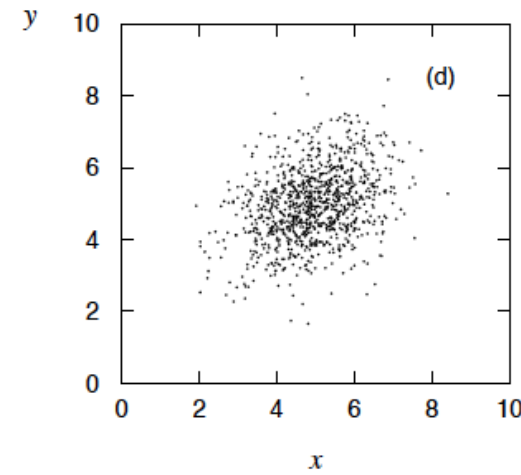- converse not always true.

# Correlation (cont.)



$\rho = 0.75$

$\rho = -0.75$

$\rho = 0.95$

$\rho = 0.25$

# Covariance matrix

Suppose we have a set of $n$ random variables, say, $x_1, ..., x_n$.

We can write the covariance of each pair as an $n$ x $n$ matrix:

$$V_{ij} = \text{cov}[x_i, x_j] = \rho_{ij}\sigma_i\sigma_j$$

$$V = \begin{pmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 & \cdots & \rho_{1n}\sigma_1\sigma_n \\ \rho_{21}\sigma_2\sigma_1 & \sigma_2^2 & \cdots & \rho_{2n}\sigma_2\sigma_n \\ \vdots & & \ddots & \vdots \\ \rho_{n1}\sigma_n\sigma_1 & \rho_{n2}\sigma_n\sigma_2 & \cdots & \sigma_n^2 \end{pmatrix}$$

Covariance matrix is:

symmetric,

diagonal = variances,

positive semi-definite:

$z^T V z \geq 0$ for all $z \in \mathbb{R}^n$

# Correlation matrix

Closely related to the covariance matrix is the $n$ x $n$ matrix of correlation coefficients:

$$\rho_{ij} = \frac{\text{cov}[x_i, x_j]}{\sigma_i \sigma_j}$$

$$\rho = \begin{pmatrix} 1 & \rho_{12} & \cdots & \rho_{1n} \\ \rho_{21} & 1 & \cdots & \rho_{2n} \\ \vdots & & \ddots & \vdots \\ \rho_{n1} & \rho_{n2} & \cdots & 1 \end{pmatrix}$$

By construction, diagonal elements are $\rho_{ii} = 1$

| Distribution/pdf | Use in HEP |
| --- | --- |
| Binomial | Branching ratio |
| Multinomial | Histogram with fixed $N$ |
| Poisson | Number of events |
| 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 $\dfrac{N!}{n!(N-n)!}$

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

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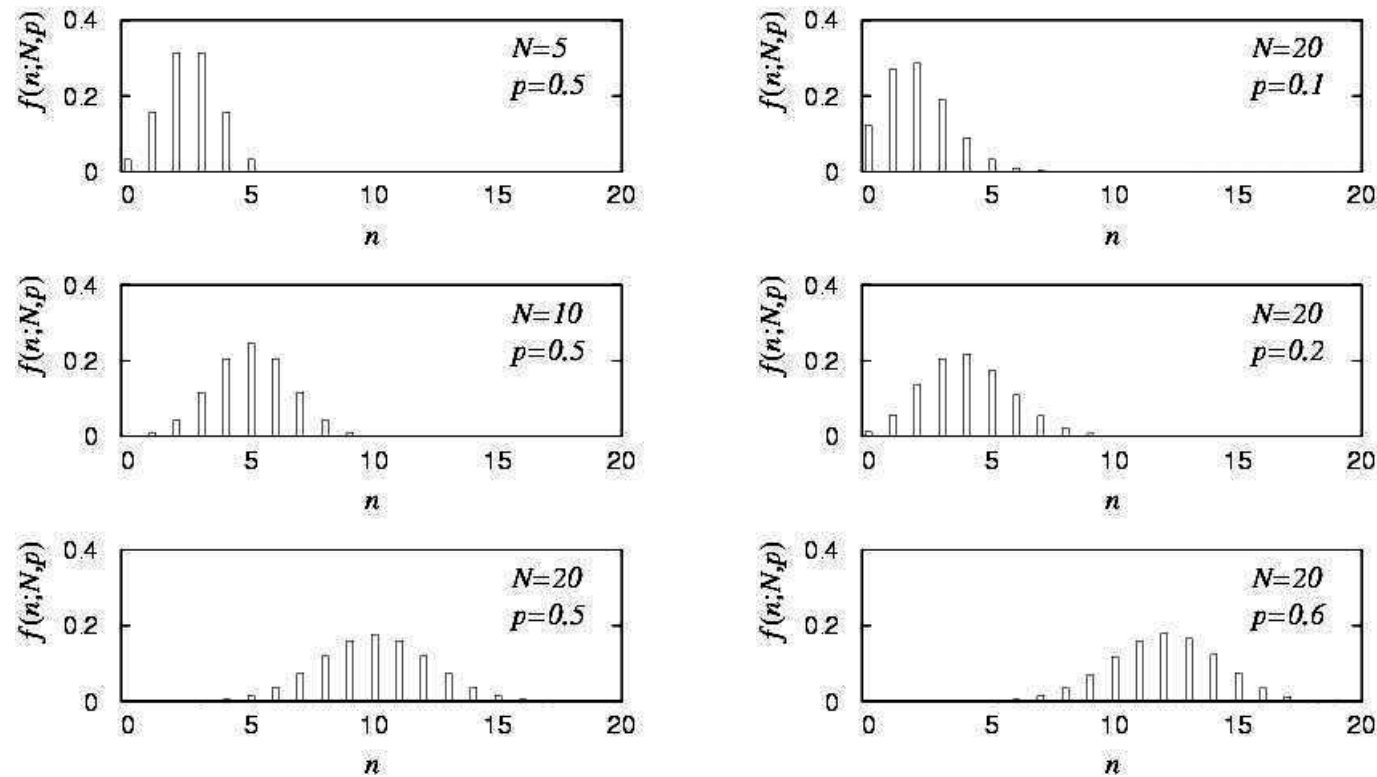
random
variable

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

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

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 $\vec{n} = (n_1, \ldots, n_m)$

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

# Multinomial distribution (2)

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

$\rightarrow$ 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.
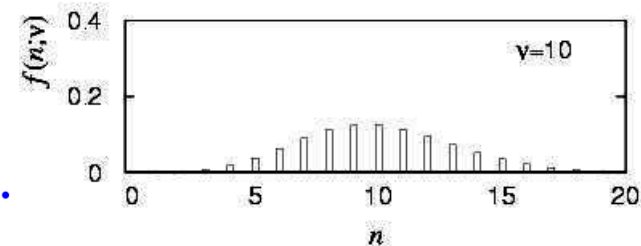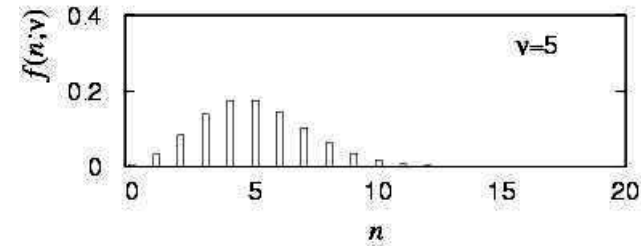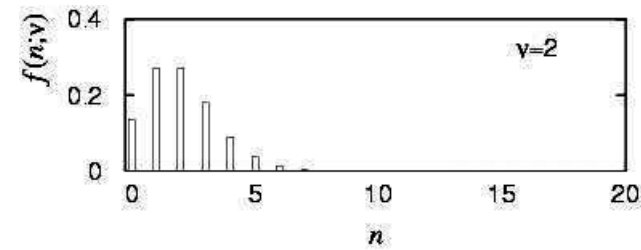
# Poisson distribution

Consider binomial $n$ in the limit

$$N \to \infty, \qquad p \to 0, \qquad E[n] = Np \to \nu .$$

$\to$ $n$ follows the Poisson distribution:

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

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

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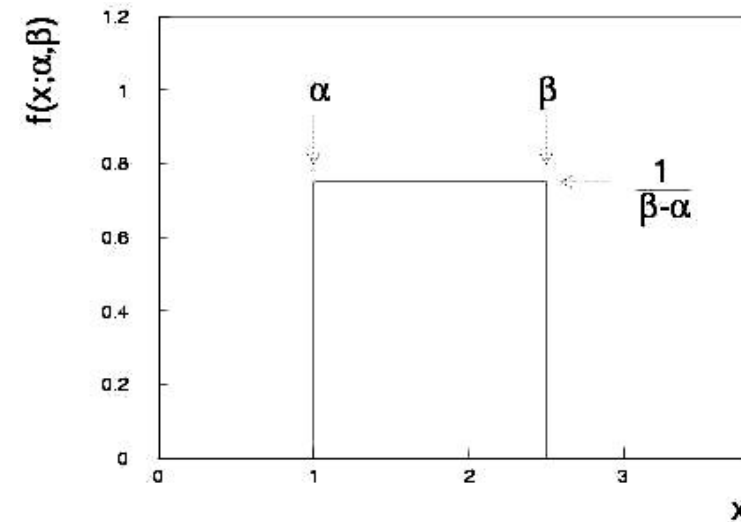
# Uniform distribution

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

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

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

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



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

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

$$E_{\min} = \frac{1}{2}E_\pi(1 - \beta), \qquad E_{\max} = \frac{1}{2}E_\pi(1 + \beta)$$
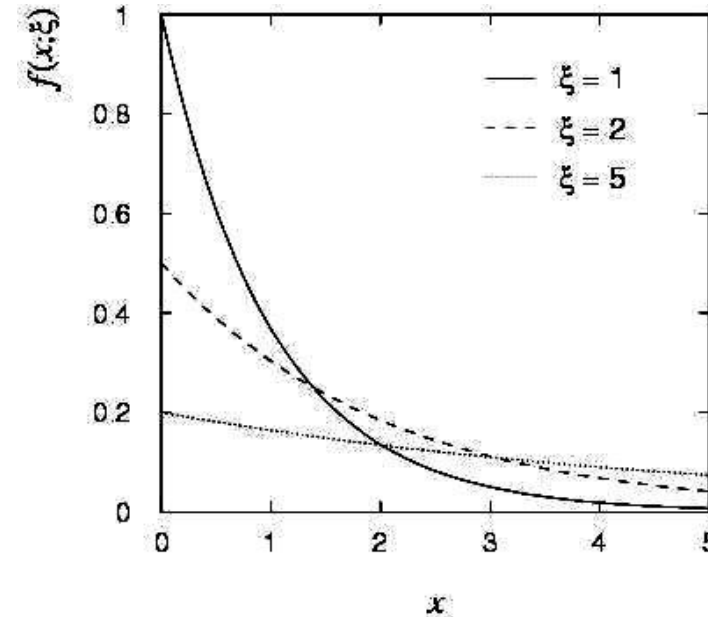
# Exponential distribution

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

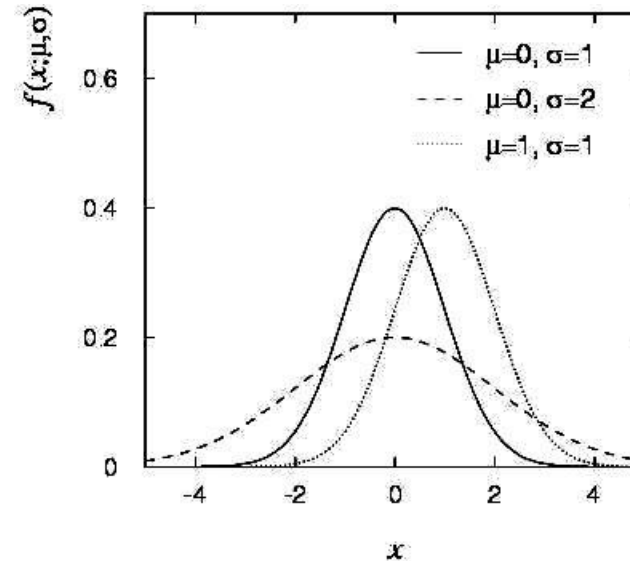$$f(t;\tau) = \frac{1}{\tau}e^{-t/\tau} \qquad (\tau = \text{mean lifetime})$$

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:

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

$E[x] = \mu$     (often $\mu$, $\sigma^2$ denote

mean, variance of any

$V[x] = \sigma^2$     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 $n \to \infty$, $y$ is a Gaussian r.v. with

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

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

# Multivariate Gaussian distribution

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

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

$\vec{x}, \vec{\mu}$ 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

$$f(x_1, x_2, ; \mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1 - \rho^2}}$$

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

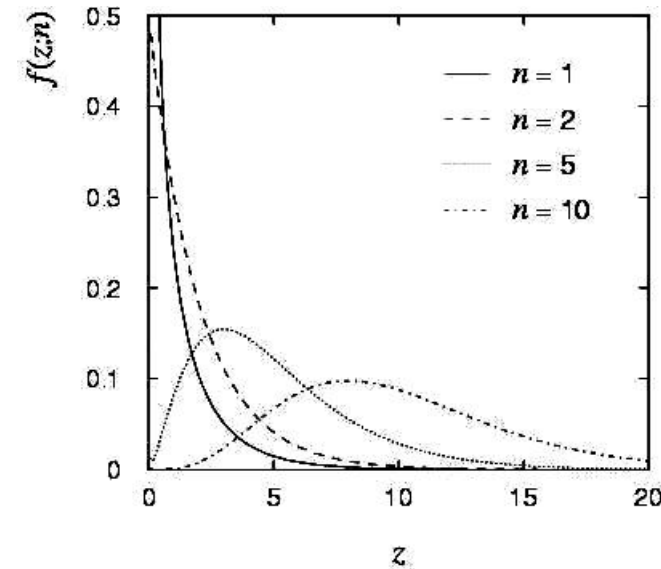where $\rho = \text{cov}[x_1, x_2]/(\sigma_1\sigma_2)$ is the correlation coefficient.

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

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

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

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

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



For independent Gaussian $x_i$, $i = 1, ..., n$, 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.
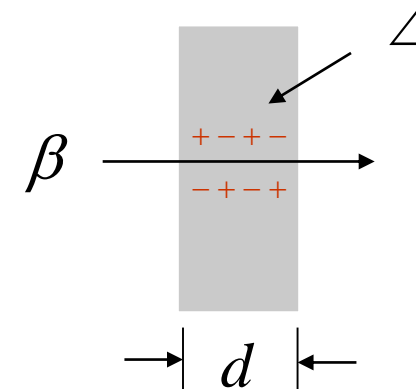
# Landau distribution

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

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

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

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

$$\xi = \frac{2\pi N_A e^4 z^2 \rho \sum Z}{m_e c^2 \sum A}\frac{d}{\beta^2} \, , \qquad \epsilon' = \frac{I^2 \exp \beta^2}{2m_e c^2 \beta^2 \gamma^2} \, .$$
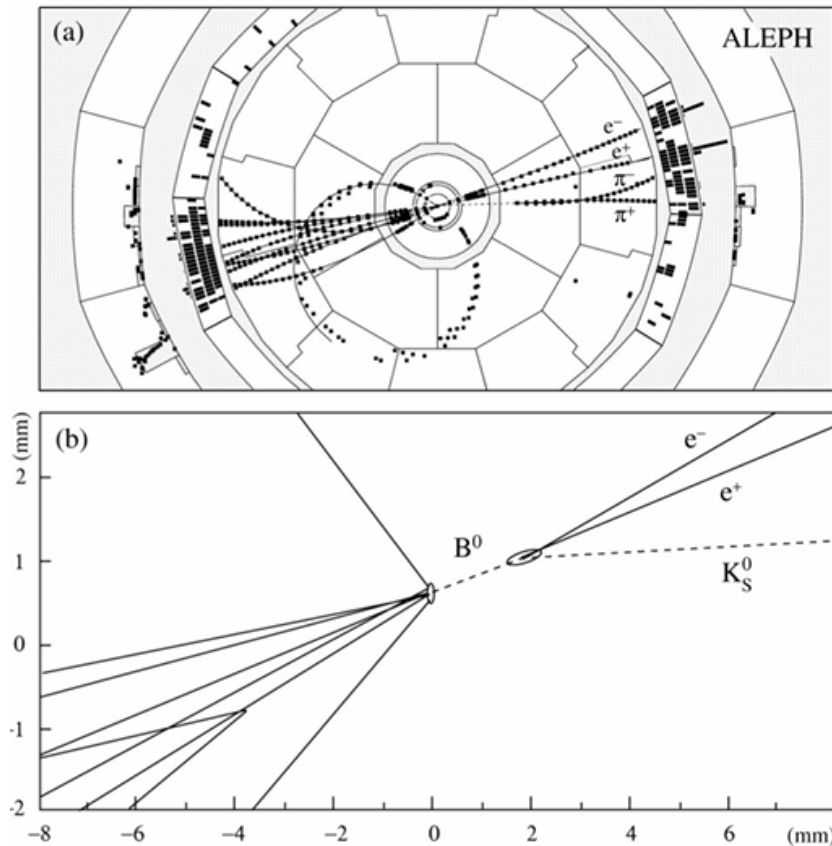
L. Landau, J. Phys. USSR **8** (1944) 201; see also
W. Allison and J. Cobb, Ann. Rev. Nucl. Part. Sci. **30** (1980) 253.

# Example: decay of an unstable particle

As an example that we'll use to illustrate several statistical methods, consider measuring the proper decay time of an unstable particle such as a B meson:

Measure flight distance $d$ and momentum $p$ of decay products of B meson with mass $m_B$.

These are related to the proper decay time $t_p$ (time in B rest frame) by

$$d = vt_{lab} = \beta c \times \gamma t_p = \frac{p_B}{m_B} t_p$$
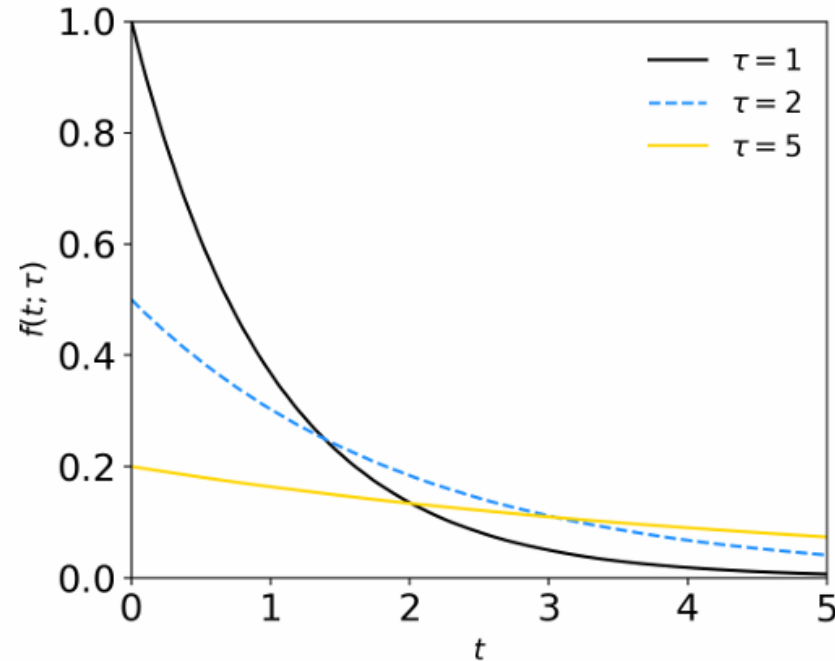
so $\quad t_p = \frac{m_B d}{p_B}$

# Exponential pdf for proper decay time

We can model $t$ as following an exponential pdf:

$$f(t; \tau) = \frac{1}{\tau} e^{-t/\tau} , \qquad t \geq 0$$

random variable    parameter



We can show (exercise) that the mean and variance of $t$ are:

$$E[t] = \int_0^\infty t f(t; \tau)\, dt = \tau \qquad\qquad V[t] = E[t^2] - (E[t])^2 = \tau^2$$

# Frequentist hypothesis tests

Suppose a measurement produces data $x$; consider a hypothesis $H_0$ we want to test and alternative $H_1$

$H_0, H_1$ specify probability for $x$: $P(x|H_0)$, $P(x|H_1)$

A test of $H_0$ is defined by specifying a critical region $w$ of the data space such that there is no more than some (small) probability $\alpha$, assuming $H_0$ is correct, to observe the data there, i.e.,

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

Need inequality if data are discrete.

$\alpha$ is called the size or significance level of the test.

If $x$ is observed in the critical region, reject $H_0$.

data space $\Omega$
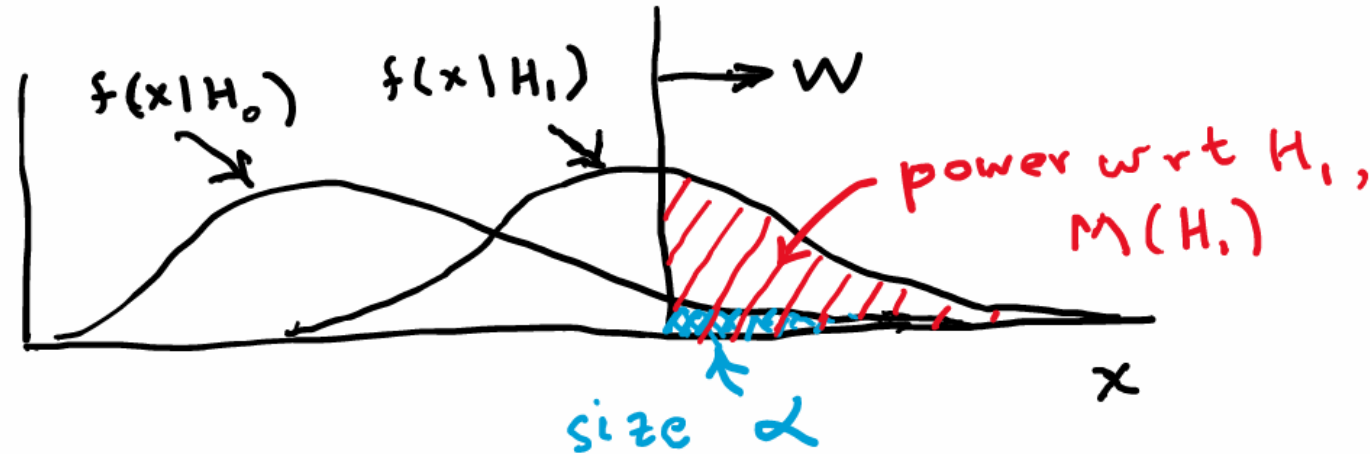
critical region $w$

# Definition of a test (2)

But in general there are an infinite number of possible critical regions that give the same size $\alpha$.

Use the alternative hypothesis $H_1$ to motivate where to place the critical region.

Roughly speaking, place the critical region where there is a low probability ($\alpha$) to be found if $H_0$ is true, but high if $H_1$ is true:

# Example of a test for classification
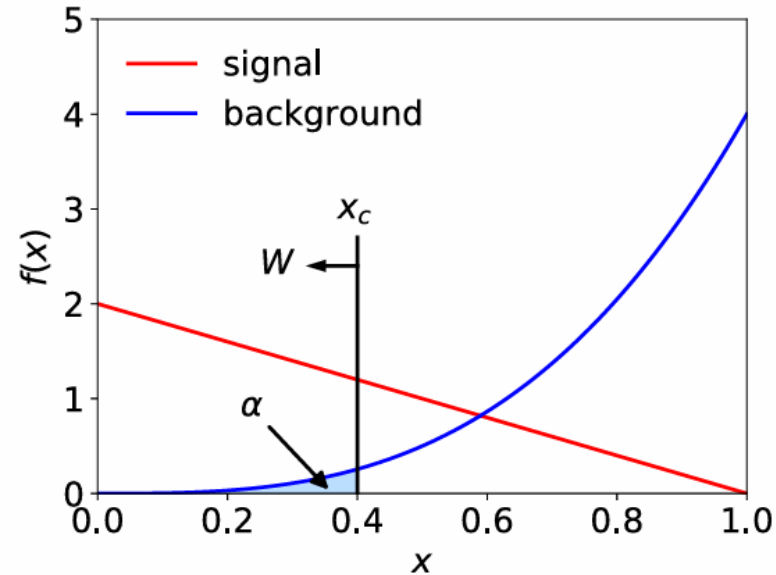
Suppose we can measure for each event a quantity $x$, where

$$f(x|s) = 2(1 - x)$$

$$f(x|b) = 4x^3$$

with $0 \leq x \leq 1$.



For each event in a mixture of signal (s) and background (b) test

$H_0$ : event is of type b

using a critical region $W$ of the form:  $W = \{x : x \leq x_c\}$, where $x_c$ is a constant that we choose to give a test with the desired size $\alpha$.

# Classification example (2)

Suppose we want $\alpha = 10^{-4}$.   Require:

$$\alpha = P(x \le x_c | b) = \int_0^{x_c} f(x|b)\, dx = \left. \frac{4x^4}{4} \right|_0^{x_c} = x_c^4$$

and therefore   $x_c = \alpha^{1/4} = 0.1$

For this test (i.e. this critical region $W$), the power with respect to the signal hypothesis (s) is

$$M = P(x \le x_c | s) = \int_0^{x_c} f(x|s)\, dx = 2x_c - x_c^2 = 0.19$$

Note:  the optimal size and power is a separate question that will depend on goals of the subsequent analysis.

# Parameter estimation

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

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

random variable   parameter

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

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

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

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

# Maximum Likelihood Estimator

## The likelihood function for i.i.d.* data

* i.i.d. = independent and identically distributed

Consider $n$ independent observations of $x$: $x_1$, ..., $x_n$, where $x$ follows $f(x; \theta)$. The joint pdf for the whole data sample is:
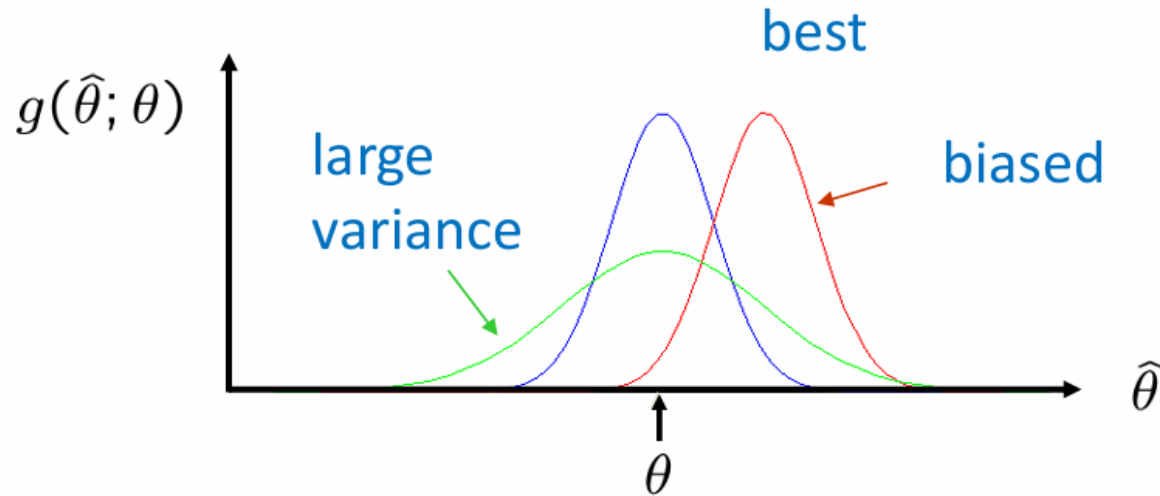
$$f(x_1, \ldots, x_n; \theta) = \prod_{i=1}^{n} f(x_i; \theta)$$

In this case the likelihood function is

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

# 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): $\quad b = E[\widehat{\theta}] - \theta$

→ average of repeated measurements should tend to true value.

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

→ small bias & variance are in general conflicting criteria

# Maximum Likelihood Estimators (MLEs)

We *define* the maximum likelihood estimators or MLEs to be the parameter values for which the likelihood is maximum.

Maximizing $L$ equivalent to maximizing $\log L$



$$\hat{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

Could have multiple maxima (take highest).

MLEs not guaranteed to have any 'optimal' properties, (but in practice they're very good).

# MLE example:  parameter of exponential pdf

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

and suppose we have i.i.d. data, $t_1, \ldots, t_n$

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

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

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

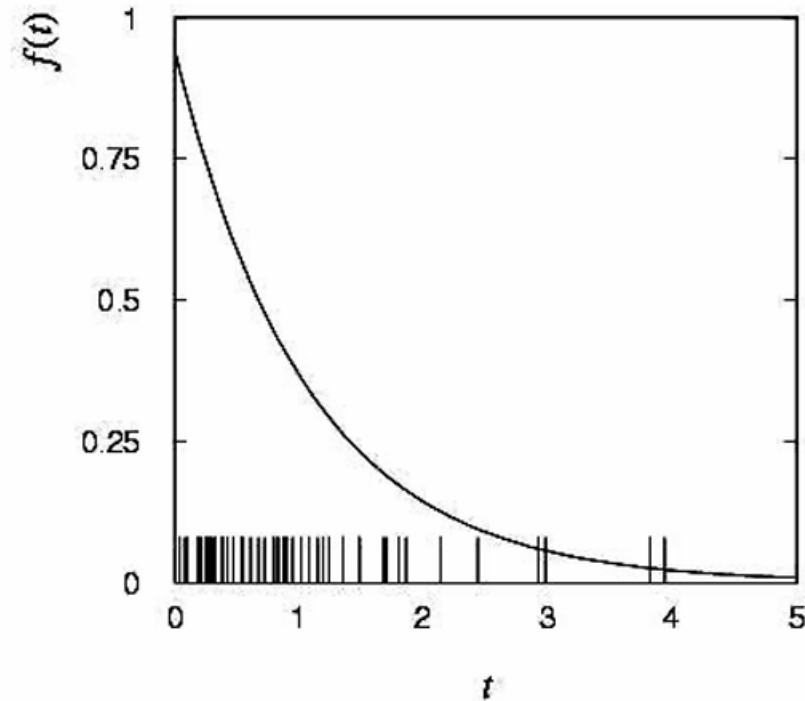# MLE example:  parameter of exponential pdf (2)

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

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

Monte Carlo test:
   generate 50  values
   using $\tau = 1$:

We find the ML estimate:

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

# MLE example: parameter of exponential pdf (3)

For the exponential distribution one has for mean, variance:

$$E[t] = \int_0^\infty t \frac{1}{\tau} e^{-t/\tau} \, dt = \tau$$

$$V[t] = \int_0^\infty (t - \tau)^2 \frac{1}{\tau} e^{-t/\tau} \, dt = \tau^2$$

For the MLE $\quad \hat{\tau} = \frac{1}{n} \sum_{i=1}^{n} t_i \quad$ we therefore find

$$E[\hat{\tau}] = E\left[\frac{1}{n} \sum_{i=1}^{n} t_i\right] = \frac{1}{n} \sum_{i=1}^{n} E[t_i] = \tau \quad \longrightarrow \quad b = E[\hat{\tau}] - \tau = 0$$

$$V[\hat{\tau}] = V\left[\frac{1}{n} \sum_{i=1}^{n} t_i\right] = \frac{1}{n^2} \sum_{i=1}^{n} V[t_i] = \frac{\tau^2}{n} \quad \longrightarrow \quad \sigma_{\hat{\tau}} = \frac{\tau}{\sqrt{n}}$$

# Variance of estimators from information inequality

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

Minimum Variance Bound (MVB)

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

$$(b = E[\hat{\theta}] - \theta)$$

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:

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

# MVB for MLE of exponential parameter

Find $\quad \mathrm{MVB} = -\left(1 + \dfrac{\partial b}{\partial \tau}\right)^2 \Bigg/ E\left[\dfrac{\partial^2 \ln L}{\partial \tau^2}\right]$

We found for the exponential parameter the MLE $\quad \hat{\tau} = \dfrac{1}{n}\displaystyle\sum_{i=1}^{n} t_i$

and we showed $b = 0$, hence $\partial b/\partial \tau = 0$.

We find $\quad \dfrac{\partial^2 \ln L}{\partial \tau^2} = \displaystyle\sum_{i=1}^{n}\left(\dfrac{1}{\tau^2} - \dfrac{2t_i}{\tau^3}\right)$

and since $E[t_i] = \tau$ for all $i$, $\quad E\left[\dfrac{\partial^2 \ln L}{\partial \tau^2}\right] = -\dfrac{n}{\tau^2}$ ,
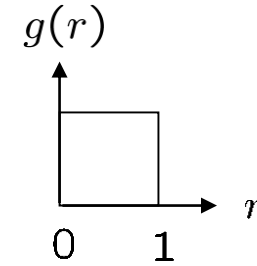
and therefore $\mathrm{MVB} = \dfrac{\tau^2}{n} = V[\hat{\tau}]$. $\qquad$ (Here MLE is "efficient").

# 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 $r_1, r_2, ..., r_m$ uniform in [0, 1].

(2) Use this to produce another sequence $x_1, x_2, ..., x_n$ 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 $\int_a^b f(x)\, dx$ .

$\rightarrow$ MC calculation = integration (at least formally)

MC generated values = 'simulated data'

$\rightarrow$ use for testing statistical procedures