Probability, statistics and data analysis

- The uncertainty in physics
- Sources of uncertaitny
- Random variables
- Probability distributions
- Error propagation
- Combination of experimental measures

The uncertainty in physics

It is generally accepted in physics that any conclusion resulting from an experimental measure is affected by a certain degree of uncertainty!



Some definitions

• Uncertainty :

parameter, associated with the result of a measurement, which characterizes the dispersion of the values that can be attributed to that measurement.

• Error :

difference between the result of a measurement and the actual value to be measured - the error is usually unknown!

• Real value

value compatible with the definition of a given quantity.

Sources of uncertainty in experimental measurements

Statistical uncertainty

 Variations in repeated observations of the quantity to be measured under (apparently) identical experimental conditions;

Systematic effects

- Incomplete definition/understanding of the quantity to be measured;
- Non-representative sampling for the quantity to be measured;
- Unawareness of the effects of environmental conditions on the measurement or imperfect measurement of environmental conditions (e.g. T, P, humidity, natural radioactivity,...);
- Parallax in the reading of analog instruments;
- Inaccurate values of measurement standards and reference materials (calibration);
- Inaccurate values of constants or other parameters obtained from sources external to the measurement(c= 3 x 10⁸ m/s,);
- Approximations and assumptions incorporated into the method / procedure

Random variables

Experimental measurements of <u>random phenomena</u> fluctuate from assay to assay.

These phenomena can be described by a random variable X, which can take a set of discrete or continuous values and which is distributed according to a frequency distribution, which associates each given value of X with a given frequency P(x). If X is continuous, then P(x) is called probability density function (p.d.f.):

 $P(x) \in [x,x+dx]=P(x)dx$

The sum of frequencies or the integral of the p.d.f. extended to the whole domain must be unity:

$$\sum_{i=1}^{n} P(x_i) = 1 \qquad \qquad \int_{x_{\min}}^{x_{\max}} P(x) dx = 1$$

Random variables

Examples:

• Throwing of dice: $x \in \{1, 2, 3, 4, 5, 6\}$



• Number of disintegrations of a radioactive source per unit of time.



- Nuclear decay is a random process;
- It is impossible to predict when one particular nucleus is going to decay;
- Only a decay probability per unit time can be assigned; this can be computed using the laws of quantum mechanics.

Defining E[x]=Sum(x P(x)), where Sum() denotes either a summation, in the case of a discrete random variable, or an integral in the case of a continuous variable, we define:

a) algebraic moment of order n : E[xⁿ]

b) central moment of order n : E[(x-E[x])ⁿ]

Defining E[x]=Sum(x P(x)), where Sum() denotes either a summation, in the case of a discrete random variable, or an integral in the case of a continuous variable, we define:

a) algebraic moment of order n : E[xⁿ]

b) central moment of order n : E[(x-E[x])ⁿ]

In particular:

• the algebraic moment of order 0, is the sum of the probabilities :

 $E[x^0] = 1$

• the algebraic moment of order 1, is the mean value (or central value) of the random variable x :

$$\mathsf{E}[\mathsf{x}] = \mu \text{ (or } \overline{x})$$

• The second central moment (of order 2), is the variance of x,

 $\mathsf{E}[(\mathsf{x}\text{-}\mathsf{E}[\mathsf{x}])^2] = \sigma^2$

(σ is the standard deviation)



Given a random variable X and two constants a and b, the following relations hold :

a) E[1] = 1;

- b) E[a+bx] = E[a] + E[bx] = a + b E[x]; in particular:
- c) E[x-E[x]] = 0 (since E[x] is a constant (= μ));
- d) $E[(x-E[x])^2] = E[x^2] (E[x])^2$

(useful relation when programming the calculation of the variance !)

N=4,n=2

Binomial distribution

N independent trials, each with only 2 possible outcomes:

heads-tails, yes-no, success-failure

p : probability for a successq=1-p : probability of failure

The probability for a sequence of N assays result in n successes and N-n failures is:

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

However, a specific sequence of n successes given N tests belongs to a set with

$$\frac{N!}{n!(N-n)!}$$

possibilities.

We therefore have:

$$P(n) = \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n}$$



$$P(n) = \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n}$$

a) Total probability is 1 :

$$\sum_{n=0}^{N} \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} = 1$$

b) Mean:

$$\bar{n} = \sum_{n=0}^{N} n \; \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} = Np$$

c) Variance:

$$\sigma^{2} = \sum_{n=0}^{N} (n-\bar{n})^{2} \frac{N!}{n!(N-n)!} p^{n} (1-p)^{N-n} = Np (1-p)$$

$$P(n) = \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n}$$

Noting that *P*(*n*) is the generic term of the <u>binomial formula</u> one has :

Binomial formula:

a)
$$\sum_{n=0}^{N} P(n) = (p + (1-p))^{N} = 1^{N} = 1$$

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k$$

$$P(n) = \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n}$$

Noting that *P*(*n*) is the generic term of the <u>binomial formula</u> one has :

a)
$$\sum_{n=0}^{N} P(n) = (p + (1-p))^{N} = 1^{N} = 1$$

b)

$$\sum_{n=0}^{N} n P(n) = \sum_{n=1}^{N} n \frac{N!}{n! (N-n)!} p^n (1-p)^{N-n} =$$

$$= pN \sum_{n=1}^{N} \frac{(N-1)!}{(n-1)! (N-n)!} p^{n-1} (1-p)^{N-n} = (\text{setting } N' = N-1 \text{ and } n' = n-1)$$

$$= pN \sum_{n'=0}^{N'} \frac{N'!}{n'! (N'-n')!} p^{n'} (1-p)^{N'-n'} = pN \sum_{n'=0}^{N'} P(n') = pN$$

(cont.)

c) Start by computing E[n(n-1)]:

$$\begin{split} E[n(n-1)] &= \\ \sum_{n=0}^{N} n(n-1) P(n) &= \sum_{n=2}^{N} n(n-1) \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} = \\ &= p^2 N(N-1) \sum_{n=2}^{N} \frac{(N-2)!}{(n-2)!(N-n)!} p^{n-2} (1-p)^{N-n} = (\text{setting } N' = N-2 \text{ and } n' = n-2) \\ &= p^2 N(N-1) \sum_{n'=0}^{N'} \frac{N'!}{n'!(N'-n')!} p^{n'} (1-p)^{N'-n'} = p^2 N(N-1) \sum_{n'=0}^{N'} P(n') = p^2 N(N-1) \end{split}$$

(cont.)

c) Start by computing E[n(n-1)]:

 $E[n(n-1)] = \sum_{n=0}^{N} n(n-1) P(n) = \sum_{n=2}^{N} n(n-1) \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} =$ = $p^2 N(N-1) \sum_{n=2}^{N} \frac{(N-2)!}{(n-2)!(N-n)!} p^{n-2} (1-p)^{N-n} = (\text{setting } N' = N-2 \text{ and } n' = n-2)$ = $p^2 N(N-1) \sum_{n'=0}^{N'} \frac{N'!}{n'!(N'-n')!} p^{n'} (1-p)^{N'-n'} = p^2 N(N-1) \sum_{n'=0}^{N'} P(n') = p^2 N(N-1)$

$$E[n(n-1)] = p^{2}N(N-1)$$

but : $E[n(n-1)] = E[n^{2} - n] = E[n^{2}] - E[n]$, so
 $E[n^{2}] = p^{2}N(N-1) + pN$
 $\sigma^{2} = E[n^{2}] - (E[n])^{2} = p^{2}N(N-1) + pN - (pN)^{2} = pN - p^{2}N = pN(1-p)$

Examples of binomial distributions

Nuclear decay

For an unstable nucleus or unstable nuclear state, the decay constant λ is the **probability per unit time** for that nucleus to decay;

In a sample of N identical nuclei, with decay constant λ , what is the probability P(n,N) of observing n decays in the unit of time ?

It is a case of **N trials** each with two possible outcomes (to decay or not to decay !), therefore described by the binomial distribution.

$$P(n,N) = \frac{N!}{n! (N-n)!} \lambda^n (1-\lambda)^{N-n}$$

The mean number of decaying nuclei per unit time is N λ



Examples of binomial distributions

Photoelectron emission in a photomultiplier



If ε is the probability of one single photon to produce one photoelectron (ε is the quantum efficiency, Q.E.) then n_{γ} photons will produce n_{e} photoelectrons, with a binomial probability :

$$P(n_e, n_{\gamma}) = \frac{n_{\gamma}!}{n_e! (n_{\gamma} - n_e)!} \varepsilon^{n_e} (1 - \varepsilon)^{n_{\gamma} - n_e}$$

The mean number of photoelectrons is $\boldsymbol{\epsilon} \, \boldsymbol{n_v}$

Poisson distribution:

It is also a discrete distribution.

The Poisson distribution describes cases in which the probability of success is very small, but the expected average value for the number of successes is constant. Usually only the average number of successes is known. The probability of each success and the number of trials are unknown.

The Poisson distribution is the limit of the binomial distribution* when N-> ∞ and p->0 with μ =Np=cte. The probability of obtaining n successes for an expected average value μ , is:

$$P(n) = \frac{\mu^n e^{-\mu}}{n!}$$



Poisson distribution

* The Poisson distribution is the limit of the binomial distribution when $N \rightarrow \infty$ and $p \rightarrow 0$ with $\mu = Np = cte$.

$$\lim_{\substack{N \to \infty \\ p \to 0}} P(n) = \lim_{\substack{N \to \infty \\ p \to 0}} \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} \dots$$

$$\lim_{N \to \infty} \frac{N!}{(N-n)!} = \lim_{N \to \infty} N(N-1)(N-2)...(N-(n-1)) \sim N^n$$

$$\lim_{N \to \infty} (1-p)^{N-n} = \lim_{k \to \infty} \sum_{k=0}^{N-n} \frac{(N-n)!}{k!(N-n-k)!} (-p)^k (1)^{N-n-k} = \sum_{k=0}^{\infty} \frac{(N-n)^k}{k!} (-p)^k \approx \sum_{k=0}^{\infty} \frac{(-Np)^k}{k!} = e^{-Np}$$

$$\lim_{\substack{N \to \infty \\ p \to 0}} P(n) = \frac{N^n}{n!} p^n e^{-Np} = \frac{(Np)^n}{n!} e^{-Np}$$

****** Binómio de Newton

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k$$

Moments of the Poisson distribution

$$P(n) = \frac{\mu^n e^{-\mu}}{n!}$$

a) Total probability is 1

$$\sum_{n=0}^{\infty} P(n) = e^{-\mu} \sum_{n=0}^{\infty} \frac{\mu^n}{n!} = e^{-\mu} \times e^{\mu} = 1$$

b) Mean:

$$\bar{n} = \sum_{n=0}^{\infty} n P(n) = \sum_{n=0}^{\infty} n \frac{\mu^n}{n!} e^{-\mu} = \mu$$

c) Variance:

$$\sigma^2 = \sum_{n=0}^{\infty} (n-\bar{n})^2 P(n) = \sum_{n=0}^{\infty} (n-\bar{n})^2 \frac{\mu^n}{n!} e^{-\mu} = \mu \qquad \text{Very important !}$$

N.B. – to prove b) and c) use the same strategy followed for the binomial distribution

Gauss distribution:

Continuous distribution, defined in $-\infty e +\infty$:

$$P(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Most instrumental errors follow a Gaussian distribution, i.e. in the measurement of lengths, times, temperatures, current voltages, etc., the results follow a normal distribution.

The width of the Gauss distribution is characterized by the standard deviation, but also the **full width at half maximum (FWHM)** can be used **: FWHM = 2.35**σ



Fig. 4.3, The Gaussian distribution for various σ . The standard deviation determines the width of the distribution



The Gauss distribution is the limit of the Poisson distrtibution for large values of $\boldsymbol{\mu}$



Gauss distribution :

Meaning of σ

Δx	∫ P(x) dx
μ±σ	68,3%
μ±2σ	95.5%
μ±3σ	99.7%



Distribution	pdf	Mean	Variance
Binomial	$P(n) = \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n}$	$\bar{n} = Np$	$\sigma^2 = Np(1-p)$
Poisson	$P(n) = \frac{\mu^n e^{-\mu}}{n!}$	$\bar{n} = \mu$	$\sigma^2 = \mu$
Gaussiana	$P(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\bar{x} = \mu$	σ^2

Estimators of the moments of a distribution

Sample: Representative set of data that allows to estimate the parameters that characterize an unknown distribution.

An estimator is **unbiased** when its value approaches the true value of the parameter it intends to estimate, as the sample size increases.

Given a population sample $x_{1,}x_{2,}...,x_{n_{1}}$ of size n, from a distribution of true central value μ , and variance σ^{2} , the corresponding unbiased estimators (sample moments) are :

• the **sample mean**, defined as the arithmetic mean of the population:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

• the **sample variance**, defined as the mean of the quadratic deviations from the mean value : \int_{n}^{∞}

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}$$

For **n->** ∞ , the sample mean tends toward the true central value of the distribution, μ , and the sample variance tends toward the true variance, σ^2 .

Geiger-Mueller counter

In this work the number of disintegrations of a radioactive source per unit of time is measured. This is a counting experiment, in which the value obtained fluctuates from test to test. This process is well described by a Poisson distribution: the probability of obtaining **n** events for an expected average value μ , is:

$$P(n) = \frac{\mu^n e^{-\mu}}{n!}$$

The best estimator of the central value is the mean and its dispersion can be characterized by the standard deviation, which, for a Poisson, are $\langle n \rangle = \mu$ and $\sigma^2 = \mu$.

Thus, to a measure of the number of counts, **n**, is associated a dispersion \sqrt{n}

$$n \pm \sqrt{n}$$

And the relative error in the number of counts is given by :

$$\frac{\sqrt{n}}{n} = \frac{1}{\sqrt{n}}$$

Error propagation

Given a quantity z calculated from variables x and y, which are directly measured: z=f(x,y), we want to determine the error associated with the determination of z, σ_z . The deviation of z from the mean (expanding in the first order around the mean values of x and y) is:

$$\Delta z = z - \bar{z} = (x - \bar{x}) \left(\frac{\partial f}{\partial x}\right)_{\bar{x}} + (y - \bar{y}) \left(\frac{\partial f}{\partial y}\right)_{\bar{y}}$$

The variance is then,

$$\sigma_{z}^{2} = \overline{\left(z - \bar{z}\right)^{2}} = \overline{\left(\left(x - \bar{x}\right)\left(\frac{\partial f}{dx}\right)_{\bar{x}} + \left(y - \bar{y}\right)\left(\frac{\partial f}{dy}\right)_{\bar{y}}\right)^{2}} = \left(\left(x - \bar{x}\right)^{2}\left(\frac{\partial f}{dx}\right)_{\bar{x}}^{2} + \left(y - \bar{y}\right)^{2}\left(\frac{\partial f}{dy}\right)_{\bar{y}}^{2} + 2\left(y - \bar{y}\right)\left(x - \bar{x}\right)\left(\frac{\partial f}{dx}\right)_{\bar{x}}\left(\frac{\partial f}{dy}\right)_{\bar{y}}\right)^{2}\right)}$$

$$\sigma_{z}^{2} = \sigma_{x}^{2}\left(\frac{\partial f}{dx}\right)_{\bar{x}}^{2} + \sigma_{y}^{2}\left(\frac{\partial f}{dy}\right)_{\bar{y}}^{2} + 2\operatorname{cov}(x, y)\left(\frac{\partial f}{dx}\right)_{\bar{x}}\left(\frac{\partial f}{dy}\right)_{\bar{y}}\right)^{2}$$
Where $\operatorname{cov}(x, y) = \overline{\left(y - \bar{y}\right)\left(x - \bar{x}\right)}$ is zero if x e y are independent variables.

If x and y are proportional $cov(x, y) = \sigma_x \sigma_y$

Examples

Defining $x \oplus y = \sqrt{x^2 + y^2}$, known as *quadratic sum*, we have the following commonly used error propagation relations:

$$z = x + y : \sigma_z = \sigma_x \oplus \sigma_y;$$

$$z = \frac{x}{y} : \sigma_z = z \times \left(\frac{\sigma_x}{x} \oplus \frac{\sigma_y}{y}\right); \text{ or } \left(\frac{\sigma_z}{z}\right) = \left(\frac{\sigma_x}{x} \oplus \frac{\sigma_y}{y}\right)$$

$$z = \ln x : \sigma_z = \frac{\sigma_x}{x}$$

Exercise : show the relations above :-)

Gamma spectroscopy

In this work, as in the work of the Geiger-Mueller counter, one is dealing with a <u>counting</u> experiment. Each MCA channel is, in fact, a counting experiment, which follows a Poisson statistic. The number of counts in each channel and the associated uncertainty are given by:

$n \pm \sqrt{n}$

On the other hand, the photopeaks corresponding to monochromatic γ s have a width (measured in number of channels) which is mainly due to instrumental uncertainties (in our case dominated by the fluctuations of the light emission in the scintillating crystal), which generally follow a Gaussian distribution. The centroid of each peak, \mathbf{x}_c , computed by the software of the acquisition program, is in fact an average over the channels, $\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N$, under the ROI^{*}, weighted by the number of counts in each channel, $\mathbf{n}_1, \mathbf{n}_2, ..., \mathbf{n}_N$.

$$x_c = \frac{1}{(n_1 + n_2 + \dots + n_N)} \sum_{i=1}^N n_i x_i$$

The dispersion for each peak is calculated from the width at half maximum as:

$$\sigma = \frac{FWHM}{2.35}$$

* ROI: Region Of Interest

Combination of experimental results *Weighted mean*

To combine N different samples, $x_1, x_2, ..., x_N$ of the same distribution, of standard deviations σ_1 , $\sigma_2, ..., \sigma_N$, in order to obtain the average of the samples one must use the weighted average and its variance:

$$\bar{x} = \frac{\sum_{i}^{N} \frac{x_i}{\sigma_i^2}}{\sum_{i}^{N} \frac{1}{\sigma_i^2}} \qquad \qquad \sigma^2(\bar{x}) = \frac{1}{\sum_{i}^{N} \frac{1}{\sigma_i^2}}$$

If $\sigma = \sigma_1 = \sigma_2 = ..., = \sigma_N$, i.e. all samples have the same dispersion, the simple average is recovered !

$$\bar{x} = \frac{\sum_{i=1}^{N} x_{i}}{N} \qquad \sigma^{2}(\bar{x}) = \frac{\sigma^{2}}{N}$$

χ^2 Distribution

Quality of fit estimator



y=ax+b or Channel= k`*Energy + b where σ is the error of the measured y (e.g. the error in the peak centroid)

Linear least squares fit

Example of application: energy calibration of the multichannel analyzer

$$\chi^{2} = \sum_{i=1}^{n} \frac{(y_{i} - ax_{i} - b)^{2}}{\sigma_{i}^{2}}$$

$$\begin{bmatrix} \frac{d\chi^{2}}{da} = 0 \Rightarrow \sum_{i=1}^{n} \frac{x_{i}(y_{i} - ax_{i} - b)}{\sigma_{i}^{2}} = 0 \\ \frac{d\chi^{2}}{db} = 0 \Rightarrow \sum_{i=1}^{n} \frac{(y_{i} - ax_{i} - b)}{\sigma_{i}^{2}} = 0 \\ \frac{d\chi^{2}}{\sigma_{i}^{2}} = 0 \Rightarrow \sum_{i=1}^{n} \frac{(y_{i} - ax_{i} - b)}{\sigma_{i}^{2}} = 0 \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} a + \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} b \\ \frac{1}{2} \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{y_{i}}{\sigma_{i}^{2}} b \\ \frac{y_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{y_{i}}{\sigma$$



Linear least squares fit



 $\begin{bmatrix} C1 = aC2 + bC3 \\ C4 = aC3 + bC5 \end{bmatrix}$

C1	$\boxed{C2}$	C3	$\begin{bmatrix} a \end{bmatrix}$
C4	<i>C</i> 3	C5	$\lfloor b \rfloor$

Solving for the straight line **parameters a and b** (Cramer's rule) :

$$a = \frac{C1C5 - C3C4}{D} \qquad \sigma_a^2 = \frac{C5}{D}$$
$$b = \frac{C2C4 - C1C3}{D} \qquad \sigma_b^2 = \frac{C2}{D}$$
$$D = C2C5 - C3^2$$

or, explicitely...

Linear least squares fit

