# Tools for Physicists: Statistics

Introduction, Probability

Wolfgang Gradl

Institut für Kernphysik

Summer semester 2023



# The scientific method: how we create 'knowledge'

#### Theory / model

- usually mathematical
- self-consistent
- simple explanations, few (arbitrary) parameters
- testable predictions / hypotheses

#### Experiment

- modify or even reject theory in case of disagrement with data
- if theory requires too many adjustments it becomes unattractive
- generate surprises

Advance of scientific knowledge is *evolutionary* process with occasional revolutions

Statistical methods are important part of this process in particular in quantitative sciences like physics



Karl Popper (1902–1994)





#### Statistics in science

Statistics is needed to:

- characterise and summarise experimental results (impractical to always deal with raw data)
- quantify uncertainty of a measurement
- assess whether two measurements of the same quantity are compatible, combine measurements
- estimate parameters of an underlying model or theory
- test hypotheses:
  - determine whether a model is compatible with data

. . . .



### Aims of this mini-series

#### Understand statistical concepts

- Ability to understand physics papers
- Know some methods / standard statistical toolbox
- Statistical inference: from data to knowledge
  - ► Should we believe a physics claim?
  - Develop intuition
  - ► Know (some) pitfalls: avoid making mistakes others have already made

#### Use tools

- Hands-on part with Python / Jupyter
- Application to your own work? You decide!



### Practical information

Three sessions:

- I. Basics, introduction, statistical distributions
- 2. Parameter estimation
- 3. Confidence intervals, hypothesis testing

About 60 minutes of lecture, then  $\geq$  30 minutes hands-on tutorial

I hope this will be useful for you, but keep in mind that there is much more to statistics than can be covered in three brief hours.





## Useful reading material

Books:

- G. Cowan, Statistical Data Analysis
- R. Barlow, Statistics: A guide to the use of statistical methods in the physical sciences
- L. Lyons, Statistics for Nuclear and Particle Physicists
- A. J. Bevan, Statistical data analysis for the physical sciences
- G. Bohm, G. Zech, Introduction to Statistics and Data Analysis for Physicists (available online)

Lectures on the web:

- G. Cowan, Royal Holloway University London: Statistical Data Analysis
- K. Reygers, U Heidelberg, Stat. Methods in Particle Physics



# Dealing with uncertainty

- Underlying theory is probabilistic (quantum mechanics / QFT) source of true randomness
- Limited knowledge about measurement process even without QM random measurement errors
- Things we could know in principle, but don't e.g. from limitations of cost, time, ...

Quantify uncertainty using tools and concepts from probability



# Mathematical definition of probability

Kolmogorov axioms:

Consider a set S (the sample space) with subsets A, B, ...(events).

Define a function on the power set of  $S, P : \mathfrak{P}(S) \mapsto [0, 1]$  with

- I.  $P(A) \ge 0$  for all  $A \subset S$
- 2. P(S) = 1
- 3.  $P(A \cup B) = P(A) + P(B)$  if  $A \cap B = \emptyset$ ,

i.e. when A and B are exclusive

From these we can derive further properties:

- $P(\bar{A}) = 1 P(A)$
- $\blacksquare P(A \cup \bar{A}) = 1$
- $\bullet P(\emptyset) = 0$
- If  $A \subset B$ , then  $P(A) \leq P(B)$
- $\bullet P(A \cup B) = P(A) + P(B) P(A \cap B)$

for the mathematically inclined: proper treatment will use measure theory





### Interpretation — intuition about probability

#### Classical definition

- Assign equal probabilities based on symmetry of problem, e.g. rolling ideal dice: P(6) = 1/6
- difficult to generalise, sounds somewhat circular

#### Frequentist: relative frequency

► A, B, ... outcomes of a repeatable experiment

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

#### Bayesian: subjective probability

► A, B, ... are hypotheses (statements that are either true or false)

P(A) = degree of belief that A is true

...all three definitions consistent with Kolmogorov's axioms



### Conditional probability, independent events

Conditional probability for two events A and B:

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

" probability of *A* given *B* " Example: rolling dice

$$P(n < 3|n \text{ even}) = \frac{P((n < 3) \cap (n \text{ even}))}{P(n \text{ even})} = \frac{1/6}{1/2} = 1/3$$

Events A and B independent  $\iff P(A \cap B) = P(A) \cdot P(B)$ A is independent of B if P(A|B) = P(A)



### Bayes' theorem

Definition of conditional probability:

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

But obviously  $P(A \cap B) = P(B \cap A)$ , so:

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

Allows to 'invert' statements about probability:

of great interest to us. Want to infer P(theory|data) from P(data|theory)

Often these two are confused, knowingly or unknowingly (advertising, political campaigns, ...)



Bayes' theorem: degree of belief in a theory





Base probability (for anyone) to have a disease D:

P(D) = 0.0001P(no D) = 0.9999

Base probability (for anyone) to have a disease D:

P(D) = 0.0001P(no D) = 0.9999

Consider a test for *D*: result is positive or negative (+ or -):

$$P(+|D) = 0.98$$
 $P(+|no D) = 0.03$  $P(-|D) = 0.02$  $P(-|no D) = 0.97$ 

Base probability (for anyone) to have a disease D:

P(D) = 0.0001P(no D) = 0.9999

Consider a test for *D*: result is positive or negative (+ or -):

P(+ D) = 0.98	P(+ no D) = 0.03
P(- D) = 0.02	P(- no D) = 0.97

Suppose your result is +; should you be worried?



Base probability (for anyone) to have a disease D:

P(D) = 0.0001P(no D) = 0.9999

Consider a test for *D*: result is positive or negative (+ or -):

P(+|D) = 0.98P(+|no D) = 0.03P(-|D) = 0.02P(-|no D) = 0.97

Suppose your result is +; should you be worried?

$$P(D|+) = \frac{P(+|D) P(D)}{P(+|D) P(D) + P(+|no D) P(no D)}$$
$$= \frac{0.98 \times 0.0001}{0.98 \times 0.0001 + 0.03 \times 0.9999} = 0.0033$$

Probability that you have disease is 0.32%, i.e. you're probably ok



### Digression: what if prevalence is (much) higher?

Assume  $100 \times$  higher prevalence in population:

P(D) = 0.01P(no D) = 0.99

Then,

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|no D)P(no D)}$$
$$= \frac{0.98 \times 0.01}{0.98 \times 0.01 + 0.03 \times 0.99} = 0.248$$

should you be worried? This can't be answered by statistics, of course ... At least take another (independent) test ...



### Criticisms — Frequentists vs. Bayesians

#### Criticisms of the frequentist interpretation

- ▶  $n \rightarrow \infty$  can never be achieved in practice. When is *n* large enough?
- Want to talk about probabilities of events that are not repeatable
  - ► *P*(rain tomorrow) but there's only one tomorrow
  - P(Universe started with a big bang) only one universe available
- P is not an intrinsic property of A, but depends on how the ensemble of possible outcomes was constructed
  - ▶ P(person I talk to is a physicist) strongly depends on whether I am at a conference or at the beach

#### Criticisms of the subjective interpretation

- 'Subjective' estimate has no place in science
- How can quantify the prior state of our knowledge?

'Bayesians address the questions everyone is interested in by using assumptions that no one believes, while Frequentists use impeccable logic to deal with an issue that is of no interest to anyone' - Louis Lyons





https://xkcd.com/1132/



# Describing data



Tools for physicists: Statistics | SoSe 2023 | 17

### Random variables and probability density functions

Random variable:

Variable whose possible values are numerical outcomes of a random phenomenon

Probability density function (pdf) of a continuous variable:

P(X found in [x, x + dx]) = p(x)dx

Normalisation:

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



## Visualisation: Histograms

#### Histogram

- representation of the frequencies of numerical outcome of a random phenomenon
- $pdf \simeq histogram$  for
  - infinite data sample
  - zero bin width
  - normalised to unit area

$$\rho(x) = \lim_{\Delta x \to 0} \frac{N(x)}{N\Delta x}$$





#### Median, mean, and mode

Arithmetic **mean** of a data sample ('sample mean'):

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

Mean of a pdf:

$$\mu \equiv \langle x \rangle \equiv \int x \rho(x) dx$$
$$\equiv \text{ expectation value } E[x]$$

#### Median:

point with 50% probability above and 50% prob. below

#### Mode:

Tools for physicists: Statistics | SoSe 2023 | 20



not necessarily the same, for skewed distributions



#### Variance, standard deviation

Variance of a **distribution** (pdf):

$$V(x) = \int dx \, p(x) \, (x - \mu)^2 = E[(x - \mu)^2]$$

Variance of a data sample

$$V(x) = \frac{1}{N} \sum_{i} (x_i - \mu)^2 = \overline{x^2} - \mu^2$$

Requires knowledge of *true* mean  $\mu$ .

Replacing  $\mu$  by sample mean  $\bar{x}$  results in underestimated variance! Instead, use this:

$$\hat{\mathcal{V}}(x) = \frac{1}{N-1} \sum_{i} (x_i - \overline{x})^2$$

Standard deviation:

 $\sigma = \sqrt{V(x)}$ 



### Multivariate distributions

Outcome of an experiment characterised by tuple  $(x_1, \ldots, x_n)$ 

 $P(A \cap B) = \frac{f(x, y)}{dx} dy$ 

with f(x, y) the 'joint pdf'

Normalisation

$$\int \cdots \int f(x_1, \ldots, x_n) \mathrm{d} x_1 \cdots \mathrm{d} x_n = 1$$

Sometimes, only the pdf of one component is wanted:

$$f_1(x_1) = \int \cdots \int f(x_1, \dots, x_n) \mathrm{d} x_2 \cdots \mathrm{d} x_n$$

 $\approx$  projection of joint pdf onto individual axis: marginalised pdf

Tools for physicists: Statistics | SoSe 2023 | 22





#### Covariance and correlation

Covariance:

$$\operatorname{cov}[x, y] = E[(x - \mu_x)(y - \mu_y)]$$

Correlation coefficient:

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

If x, y independent: pdf factorises, *i.e.*  $f(x, y) = f_x(x) f_y(y)$ , and covariance becomes

$$E[(x - \mu_x)(y - \mu_y)] = \int (x - \mu_x) f_x(x) dx \int (y - \mu_y) f_y(y) dy = 0$$

Note: converse not necessarily true

### Covariance and correlation



Same (linear) correlation coefficient, but very different 2D shapes!

✓ import pandas as pd …

d	ataset = ataset	• pd.rea
~ (	0.7s	
	x	у
0	55.3846	97.1795
1	51.5385	96.0256
2	46.1538	94.4872
3	42.8205	91.4103
4	40.7692	88.3333
37	39.4872	25.3846
38	91.2821	41.5385
39	50.0000	95.7692
10	47.9487	95.0000
41	44.1026	92.6923

· · ·

142 rows × 2 columns



	da	taset.desc	ribe()
[12]	✓ 0.7	s	
		×	v
		^	,
	count	142.000000	142.000000
	mean	54.263273	47.832253
	std	16.765142	26.935403
	min	22 307700	2 949700
		22.307700	2.740700
	25%	44.102600	25.288450
	50%	53.333300	46.025600
	75%	64.743600	68.525675
	max	98.205100	99,487200











Tools for physicists: https://www autodesk.com/research/publications/same-stats-different-graphs



#### Linear combinations of random variables

Consider two random variables x and y with known covariance cov[x, y]

For uncorrelated variables, simply add variances.

How about combination of *N* independent measurements (estimates) of a quantity,  $x_i \pm \sigma$ , all drawn from the same underlying distribution?

$$\bar{x} = \frac{1}{N} \sum x_i \text{ best estimate}$$
$$V[N\bar{x}] = N^2 \sigma$$
$$\sigma_{\bar{x}} = \frac{1}{\sqrt{N}} \sigma$$



### Combination of measurements: weighted mean

Suppose we have *N* independent measurements of the same quantity, but each with a different uncertainty:  $x_i \pm \delta_i$ Weighted sum:

$$x = w_1 x_1 + w_2 x_2$$
  
$$\delta^2 = w_1^2 \delta_1^2 + w_2^2 \delta_2^2$$

Determine weights  $w_1, w_2$  under constraint  $w_1 + w_2 = 1$  such that  $\delta^2$  is minimised:

$$w_i = \frac{1/\delta_i^2}{1/\delta_1^2 + 1/\delta_2^2}$$

If original raw data of the two measurements are available, can improve this estimate by combining raw data

alternatively, use log-likelihood curves to combine measurements



### Correlation $\neq$ causation



Correlation coefficient: 0.791

significant correlation (p < 0.0001)

0.4 kg/year/capita to produce one additional Nobel laureate

improved cognitive function associated with regular intake of dietary flavonoids?



# Some important distributions



### Uniform distribution

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

Properties:

$$E[x] = \frac{1}{2}(a+b)$$
$$V[x] = \frac{1}{12}(a+b)^{2}$$

Example:

• Strip detector: resolution for one-strip clusters: pitch  $/\sqrt{12}$ 



#### Gaussian

V

#### A.k.a. normal distribution

$$g(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$
  
Mean:  $E[x] = \mu$   
Variance:  $V[x] = \sigma^2$ 

1.0  $^{2}=02$ 0.8 = 5.0  $\mu = -2$ ,  $\sigma^2 = 0.5$ .  $\oint_{0.4}^{0.6} \psi_{1,0,2}(x)$ 0.2 0.0 -5 -4 -3 -2 0 2 3 х

Standard normal distribution:  $\mu = 0, \sigma = 1$ Cumulative distribution related to error function

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-\frac{z^2}{2}} dz = \frac{1}{2} \left[ \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) + 1 \right]$$

In Python: scipy.stats.norm(loc, scale)



#### *p*-value

Probability for a Gaussian distribution corresponding to  $[\mu - Z\sigma, \mu + Z\sigma]$ :

$$P(Z\sigma) = \frac{1}{\sqrt{2\pi}} \int_{-Z}^{+Z} e^{-\frac{x^2}{2}} = \Phi(Z) - \Phi(-Z) = \operatorname{erf}\left(\frac{Z}{\sqrt{2}}\right)$$

68.27% of area within  $\pm 1\sigma$ 95.45% of area within  $\pm 2\sigma$ 99.73% of area within  $\pm 3\sigma$ 

#### *p*-value:

probability that random process (fluctuation) produces a measurement at least this far from the true mean

p-value := 1 –  $P(Z\sigma)$ 

Available in ROOT: TMath::Prob(Z\*Z) and Python: 2\*stats.norm.sf(Z) 90% of area within  $\pm 1.645\sigma$ 95% of area within  $\pm 1.960\sigma$ 99% of area within  $\pm 2.576\sigma$ 

Deviation	p-value (%)		
$1\sigma$	31.73		
$2\sigma$	4.55		
$3\sigma$	0.270		
$4\sigma$	0.006 33		
$5\sigma$	0.000 057 3		



### Why are Gaussians so useful?

Central limit theorem: sum of *n* random variables approaches Gaussian distribution, for large *n* True, if fluctuation of sum is not dominated by the fluctuation of one (or a few) terms

- Good example: velocity component  $v_x$  of air molecules
- So-so example: total deflection due to multiple Coulomb scattering.
   Rare large angle deflections give non-Gaussian tail
- Bad example: energy loss of charged particles traversing thin gas layer. Rare collisions make up large fraction of energy loss → Landau PDF

See practical part of today's lecture



#### Binomial distribution

N independent experiments

- Outcome of each is either 'success' or 'failure'
- Probability for success is p

$$f(k; N, p) = {\binom{N}{k}} p^k (1-p)^{N-k} \qquad E[k] = Np \qquad V[k] = Np(1-p)$$

 $\binom{N}{k} = \frac{N!}{k!(N-k)!}$ 

binomial coefficient: number of permutations to have k successes in  ${\cal N}$  tries

Use binomial distribution to model processes with two outcomes Example: detection efficiency = #(particles seen by detector) / #(all particles passing detector)

In the limit  $N \to \infty, p \to 0, Np = \nu = \text{const}$ , binomial distribution can be approximated by a Poisson distribution



### Poisson distribution

$$p(k;\nu) = \frac{\nu^k}{k!} e^{-\nu}$$

$$E[k] = \nu; \quad V[k] = \nu$$

Properties:

- If  $n_1$ ,  $n_2$  follow Poisson distribution, then also  $n_1 + n_2$
- $\blacksquare$  Can be approximated by Gaussian for large  $\nu$

Examples:

- Clicks of a Geiger counter in a given time interval
- Cars arriving at a traffic light in one minute







### Poisson distribution

$$p(k;\nu) = \frac{\nu^k}{k!} e^{-\nu}$$

$$E[k] = \nu; \quad V[k] = \nu$$

Properties:

- If  $n_1$ ,  $n_2$  follow Poisson distribution, then also  $n_1 + n_2$
- $\blacksquare$  Can be approximated by Gaussian for large  $\nu$

Examples:

- Clicks of a Geiger counter in a given time interval
- Cars arriving at a traffic light in one minute

probability of *k* events occurring in fixed interval of time if events ...

- ... occur with constant rate
- ... independently of time since last event



### Poisson distribution

$$p(k;\nu) = \frac{\nu^k}{k!} e^{-\nu}$$

$$E[k] = \nu; \quad V[k] = \nu$$

Properties:

- If  $n_1$ ,  $n_2$  follow Poisson distribution, then also  $n_1 + n_2$
- $\blacksquare$  Can be approximated by Gaussian for large  $\nu$

Examples:

- Clicks of a Geiger counter in a given time interval
- Cars arriving at a traffic light in one minute

#### Rare events:

 Number of Prussian cavalrymen killed by horse-kicks

Observe 10 army corps over 20 years:

122 deaths due to horse kicks,

therefore on average 0.61 deaths / (corps  $\times$  year)

Number of deaths	Actual number	Poisson
in 1 corps in 1 year	of such cases	prediction
0	109	108.7
1	65	66.3
2	22	20.2
3	3	4.1
4	1	0.6



## Exponential distribution

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

$$E[x] = \xi; \quad V[x] = \xi^2$$

Example:

Decay time of an unstable particle at rest

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

 $\tau$  = mean lifetime



Lack of memory (unique to exponential):  $f(t - t_0 | t \ge t_0) = f(t)$ Probability for an unstable nucleus to decay in the next minute is independent of whether the nucleus was just created or has already existed for a million years.



# $\chi^2$ distribution

 $x_1, \ldots, x_n$  be *n* independent standard normal ( $\mu = 0, \sigma = 1$ ) random variables. Then the sum of their squares

$$z = \sum_{i=1}^{n} x_i^2 = \sum_{i} \frac{(x' - \mu')^2}{\sigma'^2}$$

follows a  $\chi^2$  distribution with *n* degrees of freedom.

$$f(z;n) = \frac{z^{n/2-1}}{2^{n/2}\Gamma(\frac{n}{2})}e^{-z/2}, \quad z \ge 0$$
$$E[z] = n, \quad V[z] = 2n$$

Quantify goodness of fit, compatibility of measurements, ...



#### Student's t distribution

Let  $x_1, \ldots, x_n$  be distributed as  $N(\mu, \sigma)$ .

Sample mean and estimate of variance:

$$\bar{x} = \frac{1}{n} \sum_{i} x_{i}, \quad \hat{\sigma}^{2} = \frac{1}{n-1} \sum_{i} (x_{i} - \bar{x})^{2}$$

Don't know true  $\mu$ , therefore have to estimate variance by  $\hat{\sigma}$ .

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

For 
$$n \to \infty$$
,  $f(t; n) \to N(t; 0, 1)$   
Applications:

 Hypothesis tests: assess statistical significance between two sample means

Set confidence intervals (more of that later)

 $\frac{x-\mu}{\hat{\sigma}/\sqrt{n}}$  not Gaussian: Student's t-distribution with n-1 d.o.f.



#### Landau distribution

Describes energy loss of a (heavy) charged particle in a thin layer of material due to ionisation tail with large energy loss due to occasional high-energy scattering, e.g. creation of delta rays

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

$$\lambda = \frac{\Delta - \Delta_0}{\xi}$$

$$\Delta: \text{ actual energy loss}$$

$$\Delta_0: \text{ location parameter}$$

$$\xi: \text{ material property}$$

$$f(\lambda) = \frac{1}{\pi} \int_0^\infty e^{-u \ln u - \lambda u} \sin(\pi u) du$$

$$f(\lambda) = \frac{1}{\pi} \int_0^\infty e^{-u \ln u - \lambda u} \sin(\pi u) du$$

$$f(\lambda) = \frac{1}{\pi} \int_0^\infty e^{-u \ln u - \lambda u} \sin(\pi u) du$$

Unpleasant: mean and variance (all moments, really) are not defined

 $\Delta_{0}$ ξ:



### Tools

Usable and useful tools (e.g. for your analysis) depend on environment / external constraints and factors

- within working group
- international collaboration
- personal preferences
- •

Don't underestimate the cost of choosing a different approach than everyone else around you!

It may be worth it, though; just be aware of the implications!



### Tools

From my own experience with data analysis in HEP experiments:

- Use a version control system, such as git
- To paraphrase Willem van der Poel's 'Zero One Infinity' rule: The only numbers you should care about are Zero, One, and Infinity If you have to do something more than once, automate!
- Corollary: interactive tools are nice, but scripts are much better 'in production', especially to produce plots

By all means explore your data using JupyterLab or other interactive tools, but then export the result as executable script

Make use of well-maintained libraries, toolkits &c for common tasks

Yes, you can write your own algorithms to perform function minimisation or matrix inversion — but should you?



#### Interactive session

https://bit.ly/3ANSWhN choose 'Tools for Physicists: Statistics environment'

