

# Notes: Practical Statistics for Physics & Astronomy

R.B. Metcalf

Alma Mater Studiorum - Università di Bologna

February 19, 2018

## Contents

<b>1</b>	<b>What is Probability?</b>	<b>2</b>
1.1	Frequentist interpretation of probability . . . . .	2
1.2	classical interpretation of probability . . . . .	2
1.3	Subjective or Bayesian interpretation of probability . . . . .	3
1.4	Quantum mechanical probability . . . . .	4
1.5	the rules of probability . . . . .	4
<b>A</b>	<b>Matrix basics</b>	<b>7</b>
<b>B</b>	<b>notation</b>	<b>8</b>
<b>C</b>	<b>Some Useful Integrals and mathematical definitions</b>	<b>8</b>

# 1 What is Probability?

## 1.1 Frequentist interpretation of probability

Imagine there is some event, instance or outcome of an experiment or observation called A. The probability of A is the fraction of times A occurs when the experiment or observation repeated in the same way or circumstances an *infinite* number of times.

$$P(A) = \lim_{N \rightarrow \infty} \frac{\text{number of trials where A is true}}{N(\text{total number of trials})} \quad (1.1)$$

This is the traditional definition of probability as formally stated by Laplace in 1774 and almost universally used for centuries despite no one ever having done anything *exactly* the same twice let alone an *infinite* number of times.

Applying this definition to any physical phenomenon requires a partitioning of the world into things that are known and fixed on each repatriation of the observation and those things that are not known and change every repetition. If nature is deterministic and an experiment could be set up *exactly* the same way in all respects than the outcome would always be the same and probability would not apply. Of course even in classical physics it is not possible to know the state of every atom and photon that might possibly influence your measurement apparatus (or brain). It is these things that change when repeating the observation.

This partitioning between known and unknown factors seems reasonable when we talk about the positions and momenta of particles in a gas or flipping a coin, but in many other common situations where probability is used it seems less well defined. Say someone tells you that there is a 30% probability that candidate A will win an election tomorrow. Of course an identical election will never be run again and was never run in the past. There are many factors, known and unknown, that could affect an election. This statement was probably based on polling data. By the above definition of probability, this means that if the election were held an infinite number of times in which the polling data were exactly the same the candidate would win a 30% of them. This seems like a completely unverifiable claim. If scientific knowledge must be reproducible to be considered true then it would seem that any such argument should be considered unscientific. And yet probability through statistics is at the foundation of all quantitative measurements.

Lets be a bit more practical. Lets say we don't need an infinite number of trials, but just a very *large number* of them. Lets say we flip a coin a *very large number* of times. If we did it say one billion times we would not expect that *exactly* 500 million times it would be heads. We would expect that roughly half, but not exactly half of the times it would be heads even if the probability of getting heads in each flip is 1/2. We might try to quantify how close the number of heads should be to 500 million, but in doing so we would need to use a probabilistic argument that would use the very concept we are trying to define.

Many statisticians and philosophers have found this definition of probability problematic. Despite this it is the definition usually used by scientists when they are forced to addressing this subject.

## 1.2 classical interpretation of probability

The classical interpretation of probability relies on identifying events that are equally likely or probable. This is often the argumentation used in classical statistical mechanics where each micro-state of the system is taken to be equally probable. If one then says that the probability of being in either of two mutually exclusive states is the sum of their probabilities and that the some of the probabilities of being in all possible states is one then you can find a numerical value for the

probability of each state. A macro-state (one with temperature equal to some value or total energy equal to some value) corresponds to many micro-states so by adding up their probabilities you can find the probability of macro states which will not necessarily be equal.

The biggest criticism of this interpretation is that it doesn't really say what probability is, it just tells you how to calculate it in a restricted class of problems. What does it mean that two states are equally probable? What does the probability of a macro-state mean? Another problem is that not all events that we commonly apply probability to can be reduced in this way to a collection of equally probable mutually exclusive events.

### 1.3 Subjective or Bayesian interpretation of probability

Thomas Bayes (1701 - 1761) (and Jacob Bernoulli 1655-1705) had a different conception of what probability is although the idea was not put on a firm theoretical foundation until the 1940's and 50's by G. Polya, R.T. Cox and E.T. Jaynes. It did not make its way into common use in science, in the form of Bayesian statistics, until relatively recently (80s and 90s).

In this school of thought, probability theory is an extension of formal logic to situations where the truth or falsehood of a proposition (e.g. "It will rain tomorrow." or "The mass of the Earth is between  $5.972 \times 10^{24}$  kg and  $5.978 \times 10^{24}$  kg.") cannot be deduced conclusively by deductive reasoning. A proposition has a probability function that depends on the evidence for and against its truth. When deductive reasoning can be applied conclusively this function is either zero (false) or one (true). In this way Boolean logic is a limiting case of probability theory. Surprisingly from just the following requirements (or *desiderata*) on the probability function of a proposition you can deduce the rules of probability and show that they are complete without ever mentioning randomness or repetition of experiments.

1. Degrees of plausibility are represented by real numbers.
2. The measure of plausibility must exhibit qualitative agreement with rationality. This means that as new information supporting the truth of a proposition is supplied, the number which represents the plausibility will increase continuously and monotonically. Also, to maintain rationality, the deductive limit must be obtained where appropriate.
3. Consistency
  - (a) *Structured consistency* : If the conclusion can be reasoned out in more than one way, every possible way must lead to the same result.
  - (b) *Propriety*: The theory must take account of all information that is relevant to the question.
  - (c) *Jaynes consistency*: Equivalent states of knowledge must be represented by equivalent plausibility assignments. For example, if  $A, B|C = B|C$ , then the plausibility of  $A, B|C$  must equal the plausibility of  $B|C$

(taken from Gregory (2006)).

These foundational proofs are very interesting, but outside the scope of this course (for those that are interested see chapter 2 of Gregory (2006) or, more comprehensively, Jaynes (2003)). One thing that is of importance here is that this definition allows one to define the probability of something that would not usually be considered a *random variable* or a repeated event. This is central to the Bayesian method of parameter estimation and model selection that we will get to later.

$A$	$B$	$A, B$	$\overline{A}, \overline{B}$	$\overline{A} \cup \overline{B}$	$A \cup B$	$\overline{A \cup B}$	$\overline{A}, \overline{B}$
F	T	F	T	T	T	F	F
F	F	F	T	T	F	T	T
T	T	T	F	F	T	F	F
T	F	F	T	T	T	F	F

Table 1: The truth table for binary logical expressions.

## 1.4 Quantum mechanical probability

The interpretation of probability as a measure certainty, or conversely ignorance,

## 1.5 the rules of probability

Suppose the  $A, B, \dots$  are events that either occur or don't occur, that is they have values true or false (or 0 and 1 if you prefer).  $P(A)$  is the probability of  $A$  occurring or being true. We can combine events in one of two ways.  $(A, B)$  means " $A$  and  $B$ ". It is true if both of them are true and false if both are false.  $(A \cup B)$  means " $A$  or  $B$ " it is true if either  $A$  or  $B$  is true. It is true if both are true.  $\overline{A}$  means "not  $A$ ". Note that  $\overline{A \cup B} = \overline{A}, \overline{B}$  and  $\overline{A}, \overline{B} = \overline{A \cup B}$  in the sense that there are no combinations of trues and falses for  $A$  and  $B$  that give different answers on either side of the equality. See table 1. In the language of Boolean algebra, they have the same truth table and are therefor equivalent statements. Their probabilities must also be the same.

$P(A, B)$  is called the **joint probability** of events  $A$  and  $B$ .  $P(A \cup B)$  is often called the **disjoint probability** of events  $A$  and  $B$ .

$P(A|B)$  is called a **conditional probability**. It means the probability of  $A$  *given* that  $B$  is true. You can imagine every probability being a conditional probability where it is "conditioned" on everything that you assume about the state of the Universe. Some of these things are assumed to be irrelevant and are left out. Some might be relevant but it is taken for granted so they are left out. The probability that a coin comes up heads does not depend on the time of day. It does depend on the assumption that it is a fair coin - no more likely to be heads than tails - although it might not always be stated. This is a simple example of a **statistical model** for the experiment, in this case flipping a coin.

The two fundamental rules of probability theory are

$$\begin{aligned} P(A, B) &= P(A)P(B|A) && \text{product rule} \\ P(A) + P(\overline{A}) &= 1 && \text{sum rule} \end{aligned} \tag{1.2}$$

These rules are actually derivable from some basic requirements or "desiderata" of how probabilities should behave, but for our purposes we can take them to be axioms. From these two rules and logic rules we can derive all the necessary properties of probability.

There are several particularly useful results that follow from these rules. From the logical requirement that  $(A, B)$  is the same as  $(B, A)$  and the product rule we get

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)} \quad \text{Bayes' theorem} \tag{1.3}$$

Applying the sum rule to  $(A \cup B)$  gives

$$P(A \cup B) = 1 - P(\overline{A \cup B}) \quad (1.4)$$

$$= 1 - P(\overline{A}, \overline{B}) \quad (1.5)$$

$$= 1 - P(\overline{A})P(\overline{B}|\overline{A}) \quad (1.6)$$

$$= 1 - P(\overline{A}) [1 - P(B|\overline{A})] \quad (1.7)$$

$$= 1 - P(\overline{A}) - P(\overline{A})P(B|\overline{A}) \quad (1.8)$$

$$= P(A) + P(\overline{A})P(B|\overline{A}) \quad (1.9)$$

$$= P(A) + P(\overline{A}, B) \quad (1.10)$$

$$= P(A) + P(B)P(\overline{A}|B) \quad (1.11)$$

$$= P(A) + P(B) [1 - P(A|B)] \quad (1.12)$$

$$= P(A) + P(B) - P(B)P(A|B) \quad (1.13)$$

$$P(A \cup B) = P(A) + P(B) - P(B, A) \quad \text{extended sum rule} \quad (1.14)$$

In words, the disjoint probability of two events is equal to the sum of their probabilities minus their joint probability.

If  $A$  and  $B$  are **independent** then the probability of  $A$  occurring does not depend on whether  $B$  has occurred so  $P(A|B) = P(A)$  through the product rule this implies  $P(B|A) = P(B)$  and

$$P(A, B) = P(A)P(B) \quad \text{independent events} \quad (1.15)$$

If two events are **mutually exclusive**, that is they cannot occur at the same time (the first flip of a coin cannot be both heads and tails) then  $P(A, B) = 0$  and the extended sum rule becomes

$$P(A \cup B) = P(A) + P(B) \quad \text{mutually exclusive events} \quad (1.16)$$

**Example:** If you roll a die once the probability of getting a 6 *or* a 5 is  $\frac{1}{6} + \frac{1}{6} = \frac{1}{3}$ . If you roll a die twice the probability of getting a 6 *and then* a 5 is  $(\frac{1}{6}) (\frac{1}{6}) = \frac{1}{36}$ . The probability of getting a 6 *and* a 5 is twice this because,  $\frac{1}{18}$ , because there are two ways of doing this, a 6 first or a 5 first.

This second case can be calculated in an alternative way. In the first roll we must get a 5 or a 6. We have calculated that the probability of this is  $\frac{1}{3}$ . Once this is done in the second roll we must get whichever number we didn't get in the first roll, one number out of 6, probability  $\frac{1}{6}$ . The probability of these two independent events happening is then given by the product rule  $(\frac{1}{3}) (\frac{1}{6}) = \frac{1}{18}$ .

Now say we have a set of observations  $\{A_i\}$  that are all mutually exclusive and together they include all possible outcome then

$$1 = P(A_1 \cup A_2 \cup A_3 \cup \dots | B) + P(\overline{A_1 \cup A_2 \cup A_3 \cup \dots} | B) \quad (1.17)$$

$$= P(A_1 | B) + P(A_2 \cup A_3 \cup \dots | B) + 0 \quad (1.18)$$

$$= P(A_1 | B) + P(A_2 | B) + P(A_3 \cup \dots | B) \quad (1.19)$$

$$= \sum_i P(A_i | B) \quad (1.20)$$

This is the origin of the normalization requirement on any probability distribution function (PDF). Note that I have put a  $B$  in as a condition on all the probabilities, but this would hold without them.

Another important result along these lines is

$$\sum_i P(B|A_i)P(A_i) = \sum_i P(B, A_i) = \sum_i P(A_i|B)P(B) = P(B) \sum_i P(A_i|B) = P(B) \quad (1.21)$$

with the same requirements on the set  $\{A_i\}$ . This is the origin of what we will later call marginalization.

## A Matrix basics

$$(\mathbf{ABC} \dots)^T = \dots \mathbf{C}^T \mathbf{B}^T \mathbf{A}^T \quad (\text{A.1})$$

$$(\mathbf{ABC} \dots)^{-1} = \dots \mathbf{C}^{-1} \mathbf{B}^{-1} \mathbf{A}^{-1} \quad (\text{A.2})$$

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \quad (\text{A.3})$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T \quad (\text{A.4})$$

Some properties of the determinant

$$|\mathbf{A}| = \prod_i \lambda_i \quad (\text{A.5})$$

$$|\mathbf{A}^{-1}| = 1/|\mathbf{A}| \quad (\text{A.6})$$

$$|\mathbf{BA}| = |\mathbf{B}||\mathbf{A}| \quad (\text{A.7})$$

$$|c\mathbf{A}| = c^n |\mathbf{A}| \quad (\text{A.8})$$

$$|\mathbf{A}^T| = |\mathbf{A}| \quad (\text{A.9})$$

Some properties of the trace

$$\text{tr}(\mathbf{A}) = \sum_i A_{ii} \quad (\text{A.10})$$

$$\text{tr}(\mathbf{A}) = \sum_i \lambda_{ii} \quad (\text{A.11})$$

$$\text{tr}(\mathbf{A}^T) = \text{tr}(\mathbf{A}) \quad (\text{A.12})$$

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}) \quad (\text{A.13})$$

$$\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}) \quad (\text{A.14})$$

$\mathbf{A}$  is an **orthogonal matrix** if and only if

$$\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I} \quad (\text{A.15})$$

An orthogonal matrix has the following properties

$$\mathbf{A}^T = \mathbf{A}^{-1} \quad (\text{A.16})$$

$$|\mathbf{A}| = \pm 1 \quad (\text{A.17})$$

The  $|\lambda_i| = 1$  for all eigenvalues and the magnitude of all eigenvectors are 1.

$\mathbf{C}$  is a **positive definite matrix** if

$$\mathbf{x}^T \mathbf{C} \mathbf{x} > 0 \quad \forall \mathbf{x}. \quad (\text{A.18})$$

It has the following properties

- all eigenvalues are positive
- $\text{tr}(\mathbf{C}) > 0$
- all diagonal elements are positive,  $\mathbf{C}_{ii} > 0, \forall i$
- $\mathbf{C}$  is invertible

The covariance matrix is always positive definite.

"A and B"	$A, B$
"A or B"	$A \cup B$
continuous random variables	$x, y, x_i, y_i$
vector of random variables	$\mathbf{x}$ or $\vec{x}$
discrete random numbers	$n, m$
parameters	$\alpha, \beta$
estimator of parameter $\alpha$	$\theta_\alpha$ or $\hat{\alpha}$
data	$D$ or $d_i$
indexes data or for multiple random numbers	$i, j$
statistical and/or theoretical model	$M$
Gaussian or Normal pdf	$\mathcal{G}(\mathbf{x} \boldsymbol{\mu}, \mathbf{C})$
$\mathbf{x}$ is normally distributed	$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})$
$x$ is $\chi^2$ distributed with $n$ degrees of freedom	$x \sim \chi_n^2$
arithmetic mean of $N$ samples	$\bar{x}_N$
likelihood of data given model	$\mathcal{L}(\mathbf{D} M_i)$ or $P(\mathbf{D} M_i)$
Bayesian evidence of data	$\mathcal{E}(\mathbf{D})$
Heaviside function, 1 when $B$ is true, 0 otherwise	$\Theta(B)$

Table 2: notation

## B notation

Notation may vary but in general I will follow the guide in table 2

## C Some Useful Integrals and mathematical definitions

$$\int_{-\infty}^{\infty} dx e^{-\frac{x^2}{2}} = \sqrt{2\pi} \quad (\text{C.1})$$

$$\begin{aligned} \int_{-\infty}^{\infty} dx e^{-(ax^2+bx+c)} &= e^{-c} \int_{-\infty}^{\infty} dx e^{-\left(\sqrt{a}x + \frac{b}{2\sqrt{a}}\right)^2 + \frac{b^2}{4a}} = e^{-c + \frac{b^2}{4a}} \int_{-\infty}^{\infty} \frac{dy}{\sqrt{a}} e^{-y^2} \\ &= \sqrt{\frac{\pi}{a}} e^{-c + \frac{b^2}{4a}} \end{aligned} \quad (\text{C.2})$$

$$\int_0^{\infty} dx x^n e^{-\frac{1}{2}Ax^2} = 2^{\frac{n-1}{2}} A^{-\frac{n+1}{2}} \Gamma\left(\frac{n+1}{2}\right) \quad n > -1 \quad (\text{C.3})$$

The Gamma function

$$\begin{aligned} \int_0^{\infty} dx x^n e^{-x^2} &= \frac{1}{2} \Gamma\left(\frac{n+1}{2}\right) \\ \Gamma(n) &= (n-1)! \quad n = 1, 2, \dots \\ \Gamma\left(\frac{1}{2} + n\right) &= \frac{(2n)!}{4^n n!} \sqrt{\pi} \quad n = 0, 1, 2, \dots \end{aligned} \quad (\text{C.4})$$

Stirling's approximation

$$\ln N! \simeq N \ln N - N \text{ for } N \gg 1 \quad (\text{C.5})$$

or more accurately

$$N! \simeq \sqrt{2\pi N} \left(\frac{N}{e}\right)^N \text{ for } N \gg 1 \quad (\text{C.6})$$



## Index

- $\chi^2$  distribution, 30
- anticorrelated, 26
- arithmetic mean, 33
- asymptotically unbiased estimator, 37
- Bayes' theorem, 7
- Bayes's factor, 50
- Bayesian inference, 40
- Bayesian model selection, 50
- Bernoulli distribution, 18
- biased, 37
- binomial coefficient, 13, 17
- binomial distribution, 11, 17
- binomial expansion, 18
- Cauchy distribution, 16
- Cauchy–Schwarz inequality, 26
- central limit theorem, 22
- central moments, 16
- completion of squares, 29
- conditional probability, 7
- correlated variables, 26
- covariance, 26
- covariance matrix, 26
- cumulative distribution function, 15
- disjoint probability, 7
- double factorial, 21
- eigendecomposition, 28
- error function, 21
- estimator, 33
- evidence, 41
- expectation value, 15
- extended sum rule, 8
- Fisher information matrix, 59
- gamma function, 31
- Gaussian distribution, 21
- hypergeometric distribution, 19
- hypothesis testing, 55
- improper prior, 49
- independent, 8, 26
- inverse noise weighting, 36
- Jeffreys prior, 49
- joint probability, 7
- kurtosis, 16
- Lagrange multipliers, 35
- likelihood, 40
- lognormal distribution, 25
- Lorentzian profile, 16
- marginalization, 47
- mean, 16
- mean deviation, 16
- median, 16, 38
- minimum variance estimator, 35
- mode, 15
- moment generating function (MGF), 17
- moments, 16
- multimodal, 15
- multinomial distribution, 27
- multivariate distribution, 26
- multivariate Gaussian, 27
- mutually exclusive, 8
- normal distribution, 21
- nuisance parameters, 47
- null hypothesis, 55
- Occam's factor, 51
- odds, 50
- orthogonal matrix, 28, 65
- permutations, 11
- Poisson distribution, 19
- positive definite matrix, 65
- posterior probability, 40
- prior, 40
- probability distribution function (PDF), 15
- probability mass function, 15
- product rule, 7
- quintiles, 39
- random variable, 15
- rank, 39

Rao-Cramer inequality, 59

skewness, 16

standard deviation, 16

standardized variable, 16

statistic, 33, 55

statistical model, 7

student's t-distribution, 32, 37

sum rule, 7

t-distribution, 32, 37

unbiased, 33

uniform prior, 48

unimodal, 15

variance, 16

weighted mean, 34

## References

Gregory P., 2006, Bayesian Logical Data Analysis for the Physical Sciences. Cambridge University Press

Jaynes E., 2003, Probability Theory - The Logic of Science