

# **Probability and Distribution Refresher**

Korbinian Strimmer

10 October 2025

# Table of contents

<b>Welcome</b>	<b>1</b>
Updates . . . . .	1
License . . . . .	1
<b>Preface</b>	<b>2</b>
About the author . . . . .	2
About the notes . . . . .	2
<b>1 Combinatorics</b>	<b>3</b>
1.1 Some basic mathematical notation . . . . .	3
1.2 Number of permutations . . . . .	4
1.3 De Moivre-Sterling approximation of the factorial . . . . .	4
1.4 Multinomial and binomial coefficient . . . . .	5
<b>2 Probability</b>	<b>6</b>
2.1 Random variables . . . . .	6
2.2 Conditional probability . . . . .	7
2.3 Probability mass and density function . . . . .	8
2.4 Cumulative distribution function . . . . .	9
2.5 Quantile function and quantiles . . . . .	10
2.6 Expectation or mean . . . . .	11
2.7 Variance . . . . .	11
2.8 Moments of a distribution . . . . .	12
2.9 Expectation of a transformed random variable . . . . .	13
2.10 Probability as expectation . . . . .	13
2.11 Jensen's inequality for the expectation . . . . .	13
2.12 Random vectors and their mean and variance . . . . .	14
2.13 Correlation matrix . . . . .	15
2.14 Parameters and families of distributions . . . . .	16
<b>3 Transformations</b>	<b>17</b>
3.1 Affine or location-scale transformation . . . . .	17
3.2 General invertible transformation . . . . .	19
3.3 Exponential tilting and exponential families . . . . .	20
3.4 Sums of random variables and convolution . . . . .	22
3.5 Loss functions and scoring rules . . . . .	23

<b>4 Univariate distributions</b>	<b>30</b>
4.1 Binomial distribution . . . . .	30
4.2 Beta distribution . . . . .	33
4.3 Normal distribution . . . . .	36
4.4 Gamma distribution . . . . .	39
4.5 Inverse gamma distribution . . . . .	45
4.6 Location-scale $t$ -distribution . . . . .	49
<b>5 Multivariate distributions</b>	<b>54</b>
5.1 Multinomial distribution . . . . .	54
5.2 Dirichlet distribution . . . . .	57
5.3 Multivariate normal distribution . . . . .	60
5.4 Wishart distribution . . . . .	64
5.5 Inverse Wishart distribution . . . . .	68
5.6 Multivariate $t$ -distribution . . . . .	70
<b>Bibliography</b>	<b>76</b>

# Welcome

The Probability and Distribution Refresher notes were written by [Korbinian Strimmer](#) from 2018–2025. This version is from 10 October 2025.

If you have any questions, comments, or corrections please get in touch! <sup>1</sup>

## Updates

The lecture notes will be updated from time to time.

The most current version is found at the web page for the

- [online version of the Probability and Distribution Refresher notes](#).

There you can also download the Probability and Distribution Refresher notes as

- [PDF in A4 format for printing](#) (double page layout), or as
- [6x9 inch PDF for use on tablets](#) (single page layout).

## License

These notes are licensed to you under [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](#).

---

<sup>1</sup>Email address: [korbinian.stimmer@manchester.ac.uk](mailto:korbinian.stimmer@manchester.ac.uk)

# Preface

## About the author

Hello! My name is Korbinian Strimmer and I am a Professor in Statistics. I am a member of the [Statistics group at the Department of Mathematics of the University of Manchester](#). You can find more information about me on [my home page](#).

## About the notes

These supplementary notes aim to provide a quick refresher of some essentials in combinatorics and probability as well as to offer an overview over selected univariate and multivariate distributions.

The notes are supporting information for a number of lecture notes of statistical courses I am or have been teaching at the [Department of Mathematics of the University of Manchester](#).

This includes the currently offered modules:

- [MATH27720 Statistics 2: Likelihood and Bayes](#) and
- [MATH38161 Multivariate Statistics](#)

as well as the retired module (not offered any more):

- [MATH20802 Statistical Methods](#).

# 1 Combinatorics

## 1.1 Some basic mathematical notation

Scalar quantity: plain font, typically lower case ( $x, \theta, n$ ), sometimes upper case ( $K, R^2$ , distribution functions  $F, P, Q$ ).

Sets: plain font, upper case ( $\Omega, \mathcal{F}$ )

Vector quantity: bold font, lower case ( $\mathbf{x}, \boldsymbol{\theta}$ ).

Matrix quantity: bold font, upper case ( $\mathbf{X}, \boldsymbol{\Sigma}$ ).

Summation:

$$\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$$

Product:

$$\prod_{i=1}^n x_i = x_1 \times x_2 \times \dots \times x_n \\ = x_1 x_2 \dots x_n$$

The multiplication sign  $\times$  between the factors is usually omitted unless it is needed for clarity.

Indicator function (in Iverson bracket notation):

$$[A] = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{if } A \text{ is not true} \end{cases}$$

## 1.2 Number of permutations

The number of possible orderings, or permutations, of  $n$  distinct items is the number of ways to put  $n$  items in  $n$  bins with exactly one item in each bin. It is given by the **factorial**

$$n! = \prod_{i=1}^n i = 1 \times 2 \times \dots \times n$$

where  $n$  is a positive integer. For  $n = 0$  the factorial is defined as

$$0! = 1$$

as there is exactly one permutation of zero objects.

The factorial can also be obtained using the **gamma function**

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

which can be viewed as continuous version of the factorial with  $\Gamma(x) = (x-1)!$  for any positive integer  $x$ .

## 1.3 De Moivre-Sterling approximation of the factorial

The factorial is frequently approximated by the following formula derived by [Abraham de Moivre \(1667–1754\)](#) and [James Stirling \(1692–1770\)](#)

$$n! \approx \sqrt{2\pi n} n^{n+\frac{1}{2}} e^{-n}$$

or equivalently on logarithmic scale

$$\log n! \approx \left(n + \frac{1}{2}\right) \log n - n + \frac{1}{2} \log(2\pi)$$

The approximation is good for small  $n$  (but fails for  $n = 0$ ) and becomes more and more accurate with increasing  $n$ . For large  $n$  the approximation can be simplified to

$$\log n! \approx n \log n - n$$

## 1.4 Multinomial and binomial coefficient

The number of possible permutation of  $n$  items of  $K$  distinct types, with  $n_1$  of type 1,  $n_2$  of type 2 and so on, equals the number of ways to put  $n$  items into  $K$  bins with  $n_1$  items in the first bin,  $n_2$  in the second and so on. It is given by the **multinomial coefficient**

$$\binom{n}{n_1, \dots, n_K} = \frac{n!}{n_1! n_2! \dots n_K!}$$

with  $\sum_{k=1}^K n_k = n$  and  $K \leq n$ . Note that it equals the number of permutation of all items divided by the number of permutations of the items in each bin (or of each type).

If all  $n_k = 1$  and hence  $K = n$  the multinomial coefficient reduces to the factorial.

If there are only two bins / types ( $K = 2$ ) the multinomial coefficients becomes the **binomial coefficient**

$$\binom{n}{n_1} = \binom{n}{n_1, n - n_1} = \frac{n!}{n_1! (n - n_1)!}$$

which counts the number of ways to choose  $n_1$  elements from a set of  $n$  elements.

For large  $n$  and  $n_k$  we can apply the De Moivre-Sterling approximation to the multinomial coefficient, yielding

$$\log \binom{n}{n_1, \dots, n_K} = -n \sum_{k=1}^K \frac{n_k}{n} \log \left( \frac{n_k}{n} \right)$$

Note this is  $n$  times the Shannon-Gibbs entropy of a categorical distribution with  $n_k/n$  as class probabilities.

## 2 Probability

### 2.1 Random variables

A **random variable** describes a random experiment. The set of all possible outcomes is the **sample space** of the random variable and is denoted by  $\Omega$ . If  $\Omega$  is countable then the random variable is **discrete**, otherwise it is **continuous**. For a discrete random variable the sample space  $\Omega = \{\omega_1, \omega_2, \dots\}$  is composed of a finite or infinite number of **elementary outcomes**  $\omega_i$ .

An event  $A \subseteq \Omega$  is a subset of  $\Omega$ . This includes as special cases the complete set  $\Omega$  ("certain event") and the empty set  $\emptyset$  ("impossible event"). The set of all possible events is denoted by  $\mathcal{F}$ . The complementary event  $A^C = \Omega \setminus A$  is the complement of the set  $A$  in the sample space  $\Omega$ . Two events  $A_1$  and  $A_2$  are mutually exclusive if the sets are disjoint with  $A_1 \cap A_2 = \emptyset$ .

For a discrete random variable, the elementary outcomes  $\omega_i$  are referred to as **elementary events**, and they are all mutually exclusive. An event  $A$  consists of a number of elementary events  $\omega_i \in A$  and the complementary event is given by  $A^C = \{\omega_i \in \Omega : \omega_i \notin A\}$ .

The **probability of an event**  $A$  is denoted by  $\Pr(A)$ . Broadly,  $\Pr(A)$  provides a measure of the size of the set  $A$  relative to the set  $\Omega$ . The probability measure  $\Pr(A)$  satisfies the three **axioms of probability**:

- 1)  $\Pr(A) \geq 0$ , probabilities are non-negative,
- 2)  $\Pr(\Omega) = 1$ , the certain event has probability 1, and
- 3)  $\Pr(A_1 \cup A_2 \cup \dots) = \sum_i \Pr(A_i)$ , the probability of countable mutually exclusive events  $A_i$  is additive.

This implies

- $\Pr(A) \leq 1$ , probability values lie within the range  $[0, 1]$ ,
- $\Pr(A^C) = 1 - \Pr(A)$ , the probability of the complement, and
- $\Pr(\emptyset) = 0$ , the impossible event has probability 0.

From the above it is evident that probability is closely linked to set theory, in particular to measure theory which serves as the theoretical foundations of probability and generalisations. For instance, if  $\Pr(\emptyset) = 0$  is assumed instead of  $\Pr(\Omega) = 1$ , this leads to the axioms for a **positive measure** (of which probability is a special case).

### 2.2 Conditional probability

Consider two events  $A$  and  $B$ , which may not be mutually exclusive. The probability of the event " $A$  and  $B$ " is given by the probability of the set intersection  $\Pr(A \cap B)$ . The probability of the event " $A$  or  $B$ " is given by the probability of the set union

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B).$$

This identity follows from the axioms.

The **conditional probability** of event  $A$  assuming event  $B$  has occurred is given by

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

Essentially, now  $B$  acts as the new sample space relative to which  $A$  is measured, restricting it from  $\Omega$ . Note that  $\Pr(A|B)$  is generally not the same as  $\Pr(B|A)$ , see Bayes' theorem below.

Importantly, it can be seen that any probability may be viewed as conditional, namely relative to  $\Omega$  as  $\Pr(A) = \Pr(A|\Omega)$ .

From the definition of conditional probability we derive the **product rule**

$$\begin{aligned} \Pr(A \cap B) &= \Pr(A|B) \Pr(B) \\ &= \Pr(B|A) \Pr(A) \end{aligned}$$

which in turn yields **Bayes' theorem**

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

This theorem is useful for changing the order of conditioning and it plays a key role in Bayesian statistics.

If  $\Pr(A \cap B) = \Pr(A) \Pr(B)$  then the two events  $A$  and  $B$  are **independent** with  $\Pr(A|B) = \Pr(A)$  and  $\Pr(B|A) = \Pr(B)$ .

## 2.3 Probability mass and density function

The **distribution** (or **law**) of a random variable  $x$  with sample space  $\Omega$  gives the probability for each value or a range of values of  $x$  according to the underlying probability measure. This is done in practise by employing probability mass functions (for discrete random variables) or probability density functions (for continuous random variables).

Here  $x$  is a scalar random variable, denoted by lower case and plain font. We also write  $x$  for the outcomes of the random variable. Thus, we use the same symbol to denote a random variable and its realisations.<sup>1</sup>

For a discrete random variable we define the event  $A = \{x : x = a\} = \{a\}$  (corresponding to a single elementary event) and get the probability

$$\Pr(A) = \Pr(x = a) = f(a)$$

directly from the **probability mass function (pmf)**. The pmf has the property that  $\sum_{x \in \Omega} f(x) = 1$  and that  $f(x) \in [0, 1]$ .

For continuous random variables employ a **probability density function (pdf)** instead. We define the event  $A = \{x : a < x \leq a + da\}$  (corresponding to an infinitesimal interval) and then assign the probability

$$\Pr(A) = \Pr(a < x \leq a + da) = f(a)da.$$

Similarly, the probability of the event  $A = \{x : a_1 < x \leq a_2\}$  is given by

$$\Pr(A) = \Pr(a_1 < x \leq a_2) = \int_{a_1}^{a_2} f(a)da.$$

The pdf has the property that  $\int_{x \in \Omega} f(x)dx = 1$  but in contrast to a pmf the density  $f(x) \geq 0$  may take on values larger than 1.

It is sometimes convenient to refer to a pdf or a pmf collectively as **probability density mass function (pdmf)** without specifying whether  $x$  is continuous or discrete.

The set of all  $x$  for which  $f(x)$  is positive is called the **support** of the pdmf.

<sup>1</sup>This notation is common in statistical machine learning and multivariate statistics, see for example Mardia, Kent, and Bibby (1979). In alternative notation, random variables are often written in upper case and the outcomes in lower case. However, that convention doesn't work well for multivariate random quantities (i.e. random vectors and random matrices) and is also ill-suited in Bayesian statistics where the uncertainty of parameters are modelled by random variables.

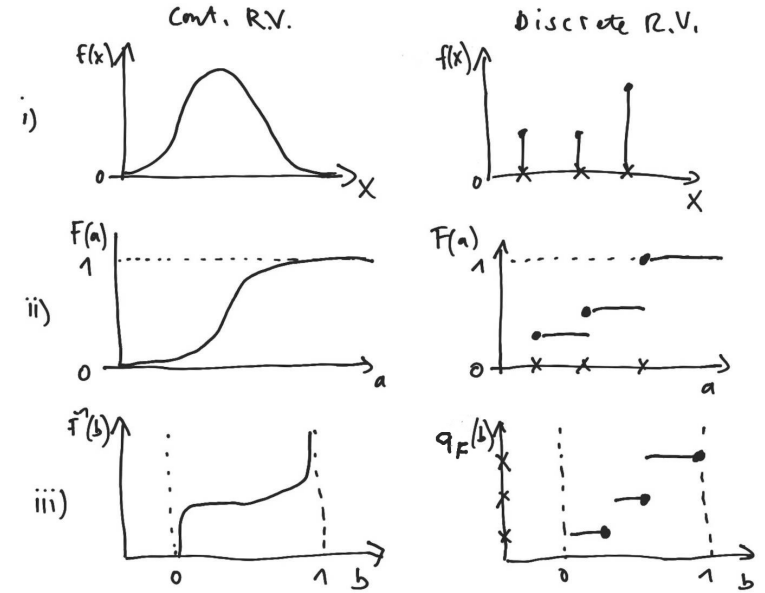


Figure 2.1: Illustration of i) pdmf, ii) distribution function and iii) quantile function for a continuous (first column) and a discrete random variable (second column).

Using the pdmf, the probability of general event  $A \subseteq \Omega$  is given by

$$\Pr(A) = \begin{cases} \sum_{x \in A} f(x) & \text{discrete case} \\ \int_{x \in A} f(x)dx & \text{continuous case} \end{cases}$$

Figure 2.1 (first row) illustrates the pdmf for a continuous and discrete random variable.

In the above we denoted the pdmf by the lower case letter  $f$  though we also often use  $p$  or  $q$ .

## 2.4 Cumulative distribution function

As alternative to the pdmf we can describe the random variable using a **cumulative distribution function (cdf)**. This requires an ordering so

that we can define the event  $A = \{x : x \leq a\}$  and compute its probability as

$$F(a) = \Pr(A) = \Pr(x \leq a) = \begin{cases} \sum_{x \in A} f(x) & \text{discrete case} \\ \int_{x \in A} f(x) dx & \text{continuous case} \end{cases}$$

The cdf is denoted by the same letter as the pdmf but in upper case (usually  $F$ ,  $P$  and  $Q$ ). By construction the cumulative distribution function is monotonically non-decreasing and its value ranges from 0 to 1. For a discrete random variable  $F(a)$  is a step function with jumps of size  $f(\omega_i)$  at the elementary outcomes  $\omega_i$ .

With the help of the cdf we can compute the probability of the event  $A = \{x : a_1 < x \leq a_2\}$  simply as

$$\Pr(A) = F(a_2) - F(a_1).$$

This works both for discrete and continuous random variables.

Figure 2.1 (second row) illustrates the distribution function for a continuous and discrete random variable.

It is common to use the same upper case letter as the cdf to name the distribution. Thus, if a random variable  $x$  has distribution  $F$  we write  $x \sim F$ , and this implies it has a pdmf  $f(x)$  and cdf  $F(x)$ .

## 2.5 Quantile function and quantiles

The **quantile function** is defined as  $q_F(b) = \min\{x : F(x) \geq b\}$ . For a continuous random variable the quantile function simplifies to  $q_F(b) = F^{-1}(b)$ , i.e. it is the ordinary inverse  $F^{-1}(b)$  of the distribution function.

Figure 2.1 (third row) illustrates the quantile function for a continuous and discrete random variable.

The quantile  $x$  of order  $b$  of the distribution  $F$  is often denoted by  $x_b = q_F(b)$ .

The 25% quantile  $x_{1/4} = x_{25\%} = q_F(1/4)$  is called the **first quartile** or **lower quartile**.

The 50% quantile  $x_{1/2} = x_{50\%} = q_F(1/2)$  is called the **second quartile** or **median**.

The 75% quantile  $x_{3/4} = x_{75\%} = q_F(3/4)$  is called the **third quartile** or **upper quartile**.

The interquartile range is the difference between the upper and lower quartiles and equals  $\text{IQR}(F) = q_F(3/4) - q_F(1/4)$ .

The quantile function is also useful for generating general random variates from uniform random variates. If  $y \sim \text{Unif}(0, 1)$  then  $x = q_F(y) \sim F$ .

## 2.6 Expectation or mean

The expected value of a random variable  $x \sim F$  is defined as the weighted average over all possible outcomes, with the weight given by the pdmf  $f(x)$ :

$$\begin{aligned} E(F) &= E(x) \\ &= \mu = \begin{cases} \sum_{x \in \Omega} f(x) x & \text{discrete case} \\ \int_{x \in \Omega} f(x) x dx & \text{continuous case} \end{cases} \end{aligned}$$

We may also write  $E_F(x)$  as a reminder that the expectation is taken with regard to the distribution  $F$ . Usually, the subscript  $F$  is left out if there are no ambiguities. A further variant is to write the expectation as  $E(F)$  to indicate that the mean is a functional of the distribution  $F$ .

Because the sum or integral may diverge, not all distributions have finite means so the mean does not always exist (in contrast to the median, or quantiles in general). For example, the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  does not have a mean for a degree of freedom in the range  $0 < \nu \leq 1$  (see Section 4.6).

## 2.7 Variance

The variance of a random variable  $x \sim F$  is the expected value of the squared deviation around the mean  $\mu = E(x)$ :

$$\begin{aligned} \text{Var}(F) &= \text{Var}(x) \\ &= E((x - \mu)^2) \\ &= E(x^2) - \mu^2 \end{aligned}$$

By construction,  $\text{Var}(x) \geq 0$ .



The notation  $\text{Var}(F)$  highlights that the variance is a functional of the distribution  $F$ . Occasionally, we write  $\text{Var}_F(x)$  indicate that the expectation is taken with regard to the distribution  $F$ .

Like the mean, the variance may diverge and hence not necessarily exists for all distribution. For example, the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  does not have a variance for the degree of freedom in the range  $0 < \nu \leq 2$  (see Section 4.6).

## 2.8 Moments of a distribution

The  $n$ -th moment of a distribution  $F$  for a random variable  $x$  is defined as follows:

$$\mu_n(F) = E(x^n)$$

Special important cases are the

- Zeroth moment:  $\mu_0(F) = E(x^0) = 1$  (since the pdmf integrates to one)
- First moment:  $\mu_1(F) = E(x^1) = E(x) = \mu$  (=the mean)
- Second moment:  $\mu_2(F) = E(x^2)$

The  $n$ -th central moment centred around the mean  $E(x) = \mu$  is given by

$$m_n(F) = E((x - \mu)^n)$$

The first few central moments are the

- Zeroth central moment:  $m_0(F) = E((x - \mu)^0) = 1$
- First central moment:  $m_1(F) = E((x - \mu)^1) = 0$
- Second central moment:  $m_2(F) = E((x - \mu)^2)$  (=the variance)

The moments of a distribution are not necessarily all finite, i.e. some moments may not exist. For example, the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  only has finite moments of degree smaller than the degree of freedom  $\nu$  (see Section 4.6).

## 2.9 Expectation of a transformed random variable

Often, one needs to find the mean of a transformed random variable. If  $x \sim F_x$  and  $y = h(x)$  with  $y \sim F_y$  then one can directly apply the above definition to obtain  $E(y) = E(F_y)$ . However, this requires knowledge of the transformed pdmf  $f_y(y)$  (see Chapter 3 for more details about variable transformations).

As an alternative, the “law of the unconscious statistician”(LOTUS) provides a convenient shortcut to compute the mean of the transformed random variable  $y = h(x)$  using only the pdmf of the original variable  $x$ :

$$E(h(x)) = \begin{cases} \sum_{x \in \Omega} f(x) h(x) & \text{discrete case} \\ \int_{x \in \Omega} f(x) h(x) dx & \text{continuous case} \end{cases}$$

Note this is not an approximation but equivalent to obtaining the mean using the transformed pdmf.

## 2.10 Probability as expectation

Probability itself can also be understood as an expectation.

For an event  $A \subseteq \Omega$  we define a corresponding indicator function  $[x \in A]$ . From LOTUS it then follows immediately that

$$\begin{aligned} E([x \in A]) &= \begin{cases} \sum_{x \in A} f(x) & \text{discrete case} \\ \int_{x \in A} f(x) dx & \text{continuous case} \end{cases} \\ &= \Pr(A) \end{aligned}$$

This relation is called the “fundamental bridge” between probability and expectation. Interestingly, one can develop the whole theory of probability from this perspective (e.g., Whittle 2000).

## 2.11 Jensen’s inequality for the expectation

If  $h(x)$  is a *convex* function then the following inequality holds:

$$E(h(\mathbf{x})) \geq h(E(\mathbf{x}))$$

Recall: a **convex** function (such as  $x^2$ ) has the shape of a “valley”.

An example of Jensen’s inequality is  $E(x^2) \geq E(x)^2$ .

## 2.12 Random vectors and their mean and variance

In addition to scalar random variables we often make use of random vectors and random matrices.<sup>2</sup>

The mean of a random vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \sim F$  is given by

$$\begin{aligned} E(F) &= E(\mathbf{x}) \\ &= \underbrace{\boldsymbol{\mu}}_{d \times 1} = (\mu_1, \dots, \mu_d)^T \end{aligned}$$

and thus is a vector of the same dimension as  $\mathbf{x}$ , where  $\mu_i = E(x_i)$  are the means of the individual components  $x_i$ .

The variance of a random vector  $\mathbf{x}$  of length  $d$ , however, is not a vector but a matrix of size  $d \times d$ . This matrix is called the **covariance matrix**:

$$\begin{aligned} \text{Var}(F) &= \text{Var}(\mathbf{x}) \\ &= \underbrace{\boldsymbol{\Sigma}}_{d \times d} = (\sigma_{ij}) = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_{dd} \end{pmatrix} \\ &= E \left( \underbrace{(\mathbf{x} - \boldsymbol{\mu})}_{d \times 1} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^T}_{1 \times d} \right) \\ &= E(\mathbf{x}\mathbf{x}^T) - \boldsymbol{\mu}\boldsymbol{\mu}^T \end{aligned}$$

The elements  $\text{Cov}(x_i, x_j) = \sigma_{ij}$  describe the covariance between the random variables  $x_i$  and  $x_j$ . The covariance matrix is symmetric, hence  $\sigma_{ij} = \sigma_{ji}$ . The diagonal elements  $\text{Cov}(x_i, x_i) = \sigma_{ii}$  correspond to the

<sup>2</sup>In our notational conventions, a scalar  $x$  is written in lower case plain font, a vector  $\mathbf{x}$  is written in lower case bold font, a matrix  $\mathbf{X}$  in upper case bold font.

individual variances  $\text{Var}(x_i) = \sigma_i^2$ . By construction, the covariance matrix  $\boldsymbol{\Sigma}$  is **positive semi-definite**, i.e. the eigenvalues of  $\boldsymbol{\Sigma}$  are all positive or equal to zero.

However, wherever possible one will aim to use models with non-singular covariance matrices, with all eigenvalues positive, so that the covariance matrix is invertible.

## 2.13 Correlation matrix

The **correlation matrix**  $\mathbf{P}$  (“upper case rho”, not “upper case p”) is the variance standardised version of the covariance matrix  $\boldsymbol{\Sigma}$ .

Specifically, denote by  $\mathbf{V}$  the diagonal matrix containing the variances

$$\mathbf{V} = \begin{pmatrix} \sigma_{11} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{dd} \end{pmatrix}$$

then the correlation matrix  $\mathbf{P}$  is given by

$$\mathbf{P} = (\rho_{ij}) = \begin{pmatrix} 1 & \dots & \rho_{1d} \\ \vdots & \ddots & \vdots \\ \rho_{d1} & \dots & 1 \end{pmatrix} = \mathbf{V}^{-1/2} \boldsymbol{\Sigma} \mathbf{V}^{-1/2}$$

Like the covariance matrix the correlation matrix is symmetric. The elements of the diagonal of  $\mathbf{P}$  are all set to 1.

Equivalently, in component notation the correlation between  $x_i$  and  $x_j$  is given by

$$\rho_{ij} = \text{Cor}(x_i, x_j) = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}\sigma_{jj}}}$$

Following from the definition above, a covariance matrix  $\boldsymbol{\Sigma}$  can be factorised into the product of standard deviations  $\mathbf{V}^{1/2}$  and the correlation matrix  $\mathbf{P}$  as follows:

$$\boldsymbol{\Sigma} = \mathbf{V}^{1/2} \mathbf{P} \mathbf{V}^{1/2}$$

## 2.14 Parameters and families of distributions

A **distribution family**  $F(\theta)$  is a collection of distributions obtained by varying a parameter  $\theta$ . Each specific value of the parameter  $\theta$  indexes one distribution in that family.

Common distribution families are usually denoted by familiar abbreviation such as  $N(\mu, \sigma^2)$  for the normal family. We also call these simply “distributions” with parameters and omit the word “family”.

If a random variable  $x$  has distribution  $F(\theta)$  we write  $x \sim F(\theta)$  or simply  $x \sim N(\mu, \sigma^2)$  in case of a named normal distribution.

The associated pdmf is written  $f(x; \theta)$  or  $f(x|\theta)$ . The conditional notation is more general because it implies the parameter  $\theta$  may have its own distribution, yielding a joint density  $f(x, \theta) = f(x|\theta)f(\theta)$ . Similarly, the corresponding cumulative distribution function is written  $F(x; \theta)$  or  $F(x|\theta)$ .

Note that parametrisations are generally not unique, as any one-to-one transformation of  $\theta$  yields an equivalent index of the same distribution family. For most commonly used distribution families there exist several standard parametrisations. We usually prefer those whose parameters that can be interpreted easily (e.g. in terms of moments) or that help to simplify calculations.

If for any pair of different parameter values  $\theta_1 \neq \theta_2$  we get distinct distributions with  $F(\theta_1) \neq F(\theta_2)$  then the distribution family  $F(\theta)$  is said to be **identifiable** by the parameter  $\theta$ .

## 3 Transformations

### 3.1 Affine or location-scale transformation

#### Transformation rule

Suppose  $x \sim F_x$  is a scalar random variable. The random variable

$$y = a + bx$$

is a **location-scale transformation** or **affine transformation** of  $x$ , where  $a$  plays the role of the **location parameter** and  $b$  is the **scale parameter**. For  $a = 0$  this is a **linear transformation**. If  $b \neq 0$  then the transformation is **invertible**, with back-transformation

$$x = (y - a)/b$$

Invertible transformations provide a one-to-one map between  $x$  and  $y$ .

For a random vector  $x \sim F_x$  of dimension  $d$  the location-scale transformation is

$$y = a + Bx$$

where  $a$  (a  $m \times 1$  vector) is the **location parameter** and  $B$  (a  $m \times d$  matrix) the **scale parameter**. For  $m = d$  (square  $B$ ) and  $\det(B) \neq 0$  the affine transformation is **invertible** with back-transformation

$$x = B^{-1}(y - a)$$

#### Density

If  $x$  is a continuous random variable with density  $f_x(x)$  and assuming an invertible transformation the density for  $y$  is given by

$$f_y(y) = |b|^{-1} f_x\left(\frac{y - a}{b}\right)$$

where  $|b|$  is the absolute value of  $b$ . Likewise, assuming an invertible transformation for a continuous random vector  $\mathbf{x}$  with density  $f_{\mathbf{x}}(\mathbf{x})$  the density for  $\mathbf{y}$  is given by

$$f_{\mathbf{y}}(\mathbf{y}) = |\det(\mathbf{B})|^{-1} f_{\mathbf{x}}(\mathbf{B}^{-1}(\mathbf{y} - \mathbf{a}))$$

where  $|\det(\mathbf{B})|$  is the absolute value of the determinant  $\det(\mathbf{B})$ .

### Moments

The transformed random variable  $y \sim F_y$  has mean

$$E(y) = a + b\mu_x$$

and variance

$$\text{Var}(y) = b^2 \sigma_x^2$$

where  $E(x) = \mu_x$  and  $\text{Var}(x) = \sigma_x^2$  are the mean and variance of the original variable  $x$ .

The mean and variance of the transformed random vector  $\mathbf{y} \sim F_y$  is

$$E(\mathbf{y}) = \mathbf{a} + \mathbf{B} \mu_x$$

and

$$\text{Var}(\mathbf{y}) = \mathbf{B} \Sigma_x \mathbf{B}^T$$

where  $E(\mathbf{x}) = \mu_x$  and  $\text{Var}(\mathbf{x}) = \Sigma_x$  are the mean and variance of the original random vector  $\mathbf{x}$ .

### Importance of affine transformations

The constants  $\mathbf{a}$  and  $\mathbf{B}$  (or  $a$  and  $b$  in the univariate case) are the parameters of the **location-scale family**  $F_y$  created from  $F_x$ . Many important distributions are location-scale families such as the normal distribution (cf. Section 4.3 and Section 5.3) and the location-scale  $t$ -distribution (Section 4.6 and Section 5.6). Furthermore, key procedures in multivariate statistics such as orthogonal transformations (including PCA) or whitening transformations (e.g. the Mahalanobis transformation) are affine transformations.

## 3.2 General invertible transformation

### Transformation rule

As above we assume  $x \sim F_x$  is a scalar random variable and  $\mathbf{x} \sim F_x$  is a random vector.

As a generalisation of invertible affine transformations we now consider general invertible transformations. For a scalar random variable we assume the transformation is specified by  $y(x) = h(x)$  and the back-transformation by  $x(y) = h^{-1}(y)$ . For a random vector we assume  $\mathbf{y}(\mathbf{x}) = \mathbf{h}(\mathbf{x})$  is invertible with back-transformation  $\mathbf{x}(\mathbf{y}) = \mathbf{h}^{-1}(\mathbf{y})$ .

### Density

If  $x$  is a continuous random variable with density  $f_x(x)$  the density of the transformed variable  $y$  can be computed exactly and is given by

$$f_y(y) = |Dx(y)| f_x(x(y))$$

where  $Dx(y)$  is the derivative of the inverse transformation  $x(y)$ .

Likewise, for a continuous random vector  $\mathbf{x}$  with density  $f_{\mathbf{x}}(\mathbf{x})$  the density for  $\mathbf{y}$  is obtained by

$$f_{\mathbf{y}}(\mathbf{y}) = |\det(D\mathbf{x}(\mathbf{y}))| f_{\mathbf{x}}(\mathbf{x}(\mathbf{y}))$$

where  $D\mathbf{x}(\mathbf{y})$  is the Jacobian matrix of the inverse transformation  $\mathbf{x}(\mathbf{y})$ .

### Moments

The mean and variance of the transformed random variable can typically only be approximated. Assume that  $E(x) = \mu_x$  and  $\text{Var}(x) = \sigma_x^2$  are the mean and variance of the original random variable  $x$  and  $E(\mathbf{x}) = \mu_x$  and  $\text{Var}(\mathbf{x}) = \Sigma_x$  are the mean and variance of the original random vector  $\mathbf{x}$ . In the **delta method** the transformation  $y(x)$  resp.  $\mathbf{y}(\mathbf{x})$  is linearised around the mean  $\mu_x$  respectively  $\mu_x$  and the mean and variance resulting from the linear transformation is reported.

Specifically, the linear approximation for the scalar-valued function is

$$y(x) \approx y(\mu_x) + Dy(\mu_x)(x - \mu_x)$$

where  $Dy(x) = y'(x)$  is the first derivative of the transformation  $y(x)$  and  $Dy(\mu_x)$  is the first derivative evaluated at the mean  $\mu_x$ , and for the vector-valued function

$$y(x) \approx y(\mu_x) + Dy(\mu_x)(x - \mu_x)$$

where  $Dy(x)$  is the Jacobian matrix (vector derivative) for the transformation  $y(x)$  and  $Dy(\mu_x)$  is the Jacobian matrix evaluated at the mean  $\mu_x$ .

In the univariate case the delta method yields as approximation for the mean and variance of the transformed random variable  $y$

$$E(y) \approx y(\mu_x)$$

and

$$\text{Var}(y) \approx (Dy(\mu_x))^2 \sigma_x^2$$

For the vector random variable  $y$  the delta method yields

$$E(y) \approx y(\mu_x)$$

and

$$\text{Var}(y) \approx Dy(\mu_x) \Sigma_x Dy(\mu_x)^T$$

Assuming  $y(x) = a + bx$ , with  $x(y) = (y - a)/b$ ,  $Dy(x) = b$  and  $Dx(y) = b^{-1}$ , recovers the univariate location-scale transformation. Likewise, assuming  $y(x) = a + Bx$ , with  $x(y) = B^{-1}(y - a)$ ,  $Dy(x) = B$  and  $Dx(y) = B^{-1}$ , recovers the multivariate location-scale transformation.

### 3.3 Exponential tilting and exponential families

Another way to change the distribution of a random variable is by **exponential tilting**.

Suppose there is a vector valued function  $t(x)$  where each component is a transformation of  $x$ , usually a simple function such the identity  $x$ , the square  $x^2$ , the logarithm  $\log(x)$  etc. These are called the **canonical statistics**. Typically, the dimension of  $t(x)$  is small.

The exponential tilt of a **base distribution**  $B$  with base function  $h(x)$  (possibly unnormalised) toward the linear combination  $\eta^T t(x)$  of the

canonical statistics  $t(x)$  and the **canonical parameters**  $\eta$  yields the distribution family  $P(\eta)$  with pdmf

$$p(x|\eta) = \underbrace{e^{\eta^T t(x)}}_{\text{exponential tilt}} h(x) / z(\eta)$$

The **normaliser** or **partition function**  $z(\eta)$  ensures that  $p(x|\eta)$  integrates to one, with

$$z(\eta) = \int_x e^{\eta^T t(x)} h(x) dx$$

In particular,  $z(\mathbf{0}) = \int_x h(x) dx$  ensures that

$$p(x|\mathbf{0}) = h(x)/z(\mathbf{0}) = b(x)$$

is a valid base pdmf. If  $h(x)$  is a pdmf then  $z(\mathbf{0}) = 1$  and  $b(x) = h(x)$ .

A distribution family  $P(\eta)$  obtained by exponential tiling is called an **exponential family**. The set of values of  $\eta$  for which  $z(\eta) < \infty$ , and hence for which  $p(x|\eta)$  is well defined, comprises the parameter space of the exponential family. Some choices of  $h(x)$  and  $t(x)$  do not yield a finite normalising factor for any  $\eta$  and hence these cannot be used to form an exponential family.

The **log-normaliser** or **log-partition function**  $a(\eta) = \log z(\eta)$  is the cumulant generating function for the canonical statistics. Its gradient yields the mean

$$E(t(x)) = \mu_t = \nabla a(\eta)$$

and the Hessian matrix the variance

$$\text{Var}(t(x)) = \Sigma_t = \nabla \nabla^T a(\eta)$$

Many common distributions are exponential families, such as the normal distribution and the Bernoulli distribution. Exponential families are central in probability and statistics. They support effective statistical learning using likelihood and Bayesian approaches, enable data reduction via minimal sufficiency and provide the basis for generalised linear models. Furthermore, exponential families often allow to generalise results established for specific cases, such as the normal distribution, to a broader domain.

See also (Wikipedia): [exponential family — table of distributions](#).

### 3.4 Sums of random variables and convolution

#### Moments

Suppose we have a sum of  $n$  independent random variables.

$$y = x_1 + x_2 + \dots + x_n$$

where each  $x_i \sim F_{x_i}$  has its own distribution and corresponding probability density mass function  $f_{x_i}(x)$ .

With  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{1}_n = (1, 1, \dots, 1)^T$  the relationship between  $y$  and  $\mathbf{x}$  can be written as affine transformation  $y = \mathbf{1}_n^T \mathbf{x}$ . Assuming  $E(x_i) = \mu_i$ ,  $\text{Var}(x_i) = \sigma_i^2$  and  $\text{Cov}(x_i, x_j) = 0$  for  $i \neq j$  the mean and variance of the random variable  $y$  equals (cf. Section 3.1)

$$E(y) = \mathbf{1}_n^T \boldsymbol{\mu} = \sum_{i=1}^n \mu_i$$

and

$$\text{Var}(y) = \mathbf{1}_n^T \text{Var}(\mathbf{x}) \mathbf{1}_n = \sum_{i=1}^n \sigma_i^2$$

Thus both the means and variance are additive (but note that for the variance this is only true because of the independence assumption).

#### Convolution

The pdmf for  $y$  is obtained by repeatedly convolving (denoted by the asterisk  $*$  operator) the pdmfs of the  $x_i$ :

$$f_y(y) = (f_{x_1} * f_{x_2} * \dots * f_{x_n})(y)$$

The **convolution** of two functions is defined as (continuous case)

$$(f_{x_1} * f_{x_2})(y) = \int_x f_{x_1}(x) f_{x_2}(y - x) dx$$

and (discrete case)

$$(f_{x_1} * f_{x_2})(y) = \sum_x f_{x_1}(x) f_{x_2}(y - x)$$

Convolution is commutative and associative so so you may convolve multiple pdmfs in any order or grouping. Furthermore, the convolution of pdmfs yields another pdmf, i.e. the resulting function integrates to one.

Many commonly used random variables can be viewed as the outcome of convolutions. For example, the sum of Bernoulli variables yields a binomial random variable and the sum of normal variables yields another normal random variable.

See also (Wikipedia): [list of convolutions of probability distributions](#).

#### Central limit theorem

The **central limit theorem**, first postulated by [Abraham de Moivre \(1667–1754\)](#) and later proved by [Pierre-Simon Laplace \(1749–1827\)](#) asserts that the distribution of the sum of  $n$  independent and identically distributed random variables with finite mean and finite variance converges in the limit of large  $n$  to a normal distribution (Section 4.3), even if the individual random variables are not themselves normal. In other words, it asserts that for large  $n$  the convolution of  $n$  identical distributions with finite first two moments converges to the normal distribution.

### 3.5 Loss functions and scoring rules

#### Loss function

A **loss or cost function**  $L(x, a)$  evaluates a prediction  $a$  (for example a parameter or a probability distribution) on the basis of an observed outcome  $x$ , and returns a numerical score.

A loss function measures, informally, the error between  $x$  and  $a$ . During optimisation the prediction  $a$  is varied and the aim is minimisation of the error (hence a loss function has *negative orientation*, smaller is better).

Adding a constant or a positive scaling factor to the loss function will not change the location of its minimum, so such loss functions are considered equivalent.

A **utility or reward function** is a loss function with a reversed sign (hence it has *positive orientation*, larger is better).

## Risk function

The **risk** of  $a$  under the distribution  $Q$  for  $x$  is defined as the expected loss

$$R_Q(a) = E_Q(L(x, a))$$

If there is no ambiguity we drop the reference to  $Q$  and write

$$R(a) = E(L(x, a))$$

The risk of  $a$  under the empirical distribution  $\hat{Q}_n$  obtained from observations  $x_1, \dots, x_n$  is the **empirical risk**

$$\hat{R}(a) = R_{\hat{Q}_n}(a) = \frac{1}{n} \sum_{i=1}^n L(x_i, a)$$

where the expectation is replaced by the sample average.

Minimising  $R(a)$  finds optimal predictions

$$a^* = \arg \min_a R(a)$$

Depending on the choice of underlying loss  $L(x, a)$  minimising the risk provides a very general optimisation-based way to identify distributional features of the distribution  $Q$  and to obtain parameter estimates.

## Scoring rules

A **scoring rule**  $S(x, P)$  is special type of loss function<sup>1</sup> that assesses the probabilistic forecast  $P$  by assigning a numerical score based on  $P$  and the observed outcome  $x$ .

The associated **risk** of  $P$  under  $Q$  is

$$R_Q(P) = E_Q(S(x, P))$$

For a **proper** scoring rule the risk  $R_Q(P)$  is minimised at  $P = Q$ , hence

$$R_Q(P) \geq R_Q(Q)$$

For a **strictly proper** scoring rule the minimum is achieved only at the true distribution  $Q$ , so equality holds exclusively for  $P = Q$ .

<sup>1</sup>As a loss function, scoring rules are negatively oriented. However, some authors consider them as utility functions with positive orientation.

A proper scoring rule induces a **divergence** between the distributions  $Q$  and  $P$ , as the difference between the risk and the minimum risk :

$$D(Q, P) = R_Q(P) - R_Q(Q) \geq 0$$

By construction, the divergence  $D(Q, P)$  is always **non-negative** and equals zero if  $P = Q$ . For a strictly proper scoring rule the divergence vanishes exclusively for  $P = Q$ .

Proper scoring rules are very useful as they allow to identify the underlying distribution and their parameters by risk minimisation or minimisation of the associated divergences.

Proper scoring rules also have a number of further useful properties. For example, various **decompositions** exist for their risk, and the divergence satisfies a **generalised Pythagorean theorem**. Furthermore, there is a **correspondence** of proper scoring rules and their associated divergences with **Bregman divergences**.

## Common loss functions

The **squared loss** or **squared error** is one of the most commonly used loss functions:

$$L(x, a) = (x - a)^2$$

The corresponding risk is the **mean squared loss** or **mean squared error** (MSE)

$$R(a) = E((x - a)^2)$$

From  $R(a) = E((x - a)^2) = E(x^2) - 2aE(x) + a^2$  it follows  $dR(a)/da = -2E(x) + 2a$  and thus that the MSE is **minimised at the mean**  $a^* = E(x)$ . The **achieved minimum risk**  $R(a^*) = \text{Var}(x)$  is the **variance**.

The **0-1 loss** function can be written as

$$L(x, a) = \begin{cases} -[x = a] & \text{discrete case} \\ -\delta(x - a) & \text{continuous case} \end{cases}$$

employing the indicator function and Dirac delta function, respectively. The corresponding risk assuming  $x \sim Q$  and pdf  $q(x)$  is

$$R_Q(a) = -q(a)$$

which is **minimised at the mode** of the pdf.

The **asymmetric loss** can be defined as

$$L(x, a; \tau) = \begin{cases} 2\tau(x - a) & \text{for } x \geq a \\ 2(1 - \tau)(a - x) & \text{for } x < a \end{cases}$$

and the corresponding risk is **minimised at the quantile**  $x_\tau$ .

For  $\tau = 1/2$  it reduces to the **absolute loss**

$$L(x, a) = |x - a|$$

whose corresponding risk is **minimised at the median**  $x_{1/2}$ .

### Logarithmic scoring rule

The most important scoring rule is the **logarithmic scoring rule** or **log-loss**

$$S(x, P) = -\log p(x)$$

The risk of  $P$  under  $Q$  based on the log-loss is the **mean log-loss** or **cross-entropy**

$$R_Q(P) = -E_Q \log p(x) = H(Q, P)$$

which is uniquely minimised for  $P = Q$ . Thus, the log-loss is **strictly proper**. Furthermore, the log-loss is notably the only **local** strictly proper scoring rule, as it solely depends on the value of the pdmf at the observed outcome  $x$ , and not on any other features of the distribution  $P$ . The minimum risk is the Shannon-Gibbs **entropy** of  $Q$ :

$$R_Q(Q) = -E_Q \log q(x) = H(Q)$$

The relationship  $H(Q, P) \geq H(Q)$ , with equality exclusively for  $P = Q$ , is known as **Gibbs' inequality**.

The divergence induced by the log-loss is the Kullback-Leibler (KL) divergence

$$D_{KL}(Q, P) = H(Q, P) - H(Q) = E_Q \log \left( \frac{q(x)}{p(x)} \right)$$

The KL divergence obeys the **data processing inequality**, i.e. applying a transformation to the underlying random variables cannot increase the KL divergence  $D_{KL}(Q, P)$  between  $Q$  and  $P$ . This property also holds for all ***f*-divergences** (of which the KL divergence is a principal example),

but is notably *not* satisfied by divergences of other proper scoring rules (and thus other Bregman divergences).

Furthermore, the KL divergence is the only divergence induced by proper scoring rules (and thus the only Bregman divergence), as well as the only *f*-divergence, that is **invariant against general coordinate transformations**. Coordinate transformations can be viewed as a special case of data processing, and for  $D_{KL}(Q, P)$  the data-processing inequality under general invertible transformations becomes an identity.

The empirical risk of a distribution family  $P(\theta)$  based on the log-loss is proportional to the log-likelihood function

$$\hat{R}(\theta) = H(\hat{Q}_n, P(\theta)) = -\frac{1}{n} \sum_{i=1}^n \log p(x_i | \theta) = -\frac{1}{n} \ell_n(\theta)$$

Minimising the empirical risk  $\hat{R}(\theta)$  is equivalent to maximising the log-likelihood function  $\ell_n(\theta)$ .

Similarly, minimising the KL divergence  $D_{KL}(\hat{Q}_n, P(\theta))$  with regard to  $\theta$  is equivalent to minimising the empirical risk  $\hat{R}(\theta)$  and hence to maximum likelihood.

### Brier or quadratic scoring rule

The **Brier scoring rule**, also known as **quadratic scoring rule**, evaluates a probabilistic categorical forecast  $P$  with corresponding class probabilities  $p_1, \dots, p_K$  given a realisation  $x$  from the categorical distribution  $Q$  with class probabilities  $q_1, \dots, q_K$ . It can be written as

$$\begin{aligned} S(x, P) &= \sum_{y=1}^K (x_y - p_y)^2 \\ &= 1 - 2 \sum_{y=1}^K x_y p_y + \sum_{y=1}^K p_y^2 \\ &= 1 - 2p_k + \sum_{y=1}^K p_y^2 \end{aligned}$$

The indicator vector  $x = (x_1, \dots, x_K)^T = (0, 0, \dots, 1, \dots, 0)^T$  contains zeros everywhere except for a single element  $x_k = 1$ . Unlike the log-loss,



the Brier score is *not local* as the pmf for  $P$  is evaluated across all  $K$  classes, not just at the realised class  $k$ .

The corresponding risk is

$$R_Q(P) = E_Q(S(x, P)) = 1 - 2 \sum_{y=1}^K q_y p_y + \sum_{y=1}^K p_y^2$$

which is uniquely minimised for  $P = Q$ . Thus, the Brier score is **strictly proper**. The minimum risk is

$$R_Q(Q) = 1 - \sum_{y=1}^K q_y^2$$

The divergence induced by the Brier score is the **squared Euclidean distance** between the two pmfs:

$$D(Q, P) = R_Q(P) - R_Q(Q) = \sum_{y=1}^K (q_y - p_y)^2$$

### Other proper scoring rules

An example of a proper, but not strictly proper, scoring rule is the squared error relative to the mean of the quoted model  $P$ :

$$S(x, P) = (x - E(P))^2$$

The corresponding risk is

$$R_Q(P) = E_Q((x - E(P))^2)$$

which is minimised for  $P = Q$  but also for any other distribution  $P$  with the same mean as  $Q$ . The minimum risk is the variance of  $Q$ :

$$R_Q(Q) = E_Q((x - E(Q))^2) = \text{Var}(Q)$$

The associated divergence is the squared distance between the two means

$$D(Q, P) = R_Q(P) - \text{Var}(Q) = (E(Q) - E(P))^2$$

which vanishes for  $P = Q$  but also for any other  $P$  with  $E(P) = E(Q)$ .

Other useful strictly proper scoring rules include

- the continuous ranked probability score (CRPS),
- the energy score, and
- the Hyvärinen scoring rule.

See also (Wikipedia): [scoring rule](#).

## 4 Univariate distributions

### 4.1 Binomial distribution

The **binomial distribution**  $\text{Bin}(n, \theta)$  is a discrete distribution counting binary outcomes.

The **Bernoulli distribution**  $\text{Ber}(\theta)$  is a special case of the binomial distribution.

#### Standard parametrisation

A binomial random variable  $x$  describes the number of successful outcomes in  $n$  identical and independent trials. We write

$$x \sim \text{Bin}(n, \theta)$$

where  $\theta \in [0, 1]$  is the probability of a positive outcome (“success”) in a single trial. Conversely,  $1 - \theta \in [0, 1]$  is the complementary probability (“failure”). The support is  $x \in \{0, 1, 2, \dots, n\}$  which notably depends on  $n$ .

The binomial distribution is often motivated by a coin tossing experiment where  $\theta$  is the probability of “head” when flipping the coin and  $x$  is the number of observed “heads” among  $n$  throws. Another common interpretation is that of an urn model where  $n$  items are distributed into two bins (Figure 4.1). Here  $\theta$  is the probability to put an item into one urn (representing “success”, “head”) and  $1 - \theta$  the probability to put it in the other urn (representing “failure”, “tail”).

The expected value is

$$E(x) = n\theta$$

and the variance is

$$\text{Var}(x) = n\theta(1 - \theta)$$

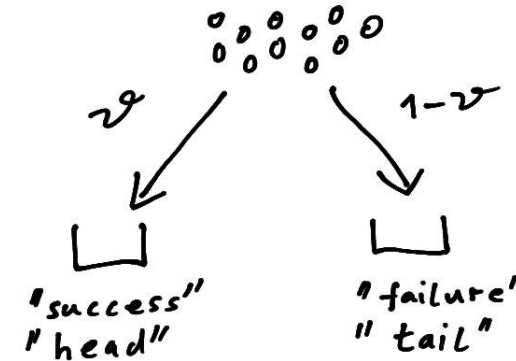


Figure 4.1: Binomial urn model.

The corresponding pmf is

$$p(x|n, \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

The binomial coefficient  $\binom{n}{x}$  in the pdf accounts for the multiplicity of ways in which we can observe  $x$  successes in  $n$  trials.

#### 💡 R code

The pmf of the binomial distribution is given by `dbinom()`, the distribution function is `pbinom()` and the quantile function is `qbinom()`. The corresponding random number generator is `rbinom()`.

#### Mean parametrisation

Instead of  $\theta$  one may also use a mean parameter  $\mu \in [0, n]$  so that

$$x \sim \text{Bin}\left(n, \theta = \frac{\mu}{n}\right)$$

The mean parameter  $\mu$  can be obtained from  $\theta$  and  $n$  by  $\mu = n\theta$ .

The mean and variance of the binomial distribution expressed in terms of  $\mu$  and  $n$  are

$$E(x) = \mu$$

and

$$\text{Var}(x) = \mu - \frac{\mu^2}{n}$$

### Special case: Bernoulli distribution

For  $n = 1$  the binomial distribution reduces to the **Bernoulli distribution**  $\text{Ber}(\theta)$ . This is the simplest of all distribution families and is named after [Jacob Bernoulli \(1655-1705\)](#) who also discovered the law of large numbers.

If a random variable  $x$  follows the Bernoulli distribution we write

$$x \sim \text{Ber}(\theta)$$

with “success” probability  $\theta \in [0, 1]$ . Conversely, the complementary “failure” probability is  $1 - \theta \in [0, 1]$ . The support is  $x \in \{0, 1\}$ . The variable  $x$  acts as an indicator variable, with “success” indicated by  $x = 1$  and “failure” indicated by  $x = 0$ .

Often the Bernoulli distribution is referred to as “coin flipping” model. Then  $\theta$  is the probability of “head” and  $1 - \theta$  the complementary probability of “tail” and  $x = 1$  corresponds to the outcome “head” and  $x = 0$  to the outcome “tail”.

The expected value is

$$\text{E}(x) = \theta$$

and the variance is

$$\text{Var}(x) = \theta(1 - \theta)$$

The pmf of  $\text{Ber}(\theta)$  is

$$p(x|\theta) = \theta^x(1 - \theta)^{1-x} = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases}$$

### Convolution property and normal approximation

The convolution of  $n$  binomial distributions, each with identical success probability  $\theta$  but possibly different number of trials  $n_i$ , yields another binomial distribution with the same parameter  $\theta$ :

$$\sum_{i=1}^n \text{Bin}(n_i, \theta) \sim \text{Bin}\left(\sum_{i=1}^n n_i, \theta\right)$$

It follows that the binomial distribution with  $n$  trials is the result of the convolution of  $n$  Bernoulli distributions:

$$\sum_{i=1}^n \text{Ber}(\theta) \sim \text{Bin}(n, \theta)$$

Thus, repeating the same Bernoulli trial  $n$  times and counting the total number of successes yields a binomial random variable.

As a consequence, following the central limit theorem (Section 3.4), for large  $n$  the binomial distribution can be well approximated by a normal distribution (Section 4.3) with the same mean and variance. This is known as the [De Moivre–Laplace theorem](#).

## 4.2 Beta distribution

The **beta distribution**  $\text{Beta}(\alpha_1, \alpha_2)$  is a continuous distribution that is useful to model proportions or probabilities for  $K = 2$  classes.

It includes the **uniform distribution** over the unit interval as a special case.

### Standard parametrisation

A beta-distributed random variable is denoted by

$$x \sim \text{Beta}(\alpha_1, \alpha_2)$$

with shape parameters  $\alpha_1 > 0$  and  $\alpha_2 > 0$ . Let  $m = \alpha_1 + \alpha_2$ . The support of  $x$  is the unit interval given by  $x \in [0, 1]$ . Thus, the beta distribution is defined over a one-dimensional space.

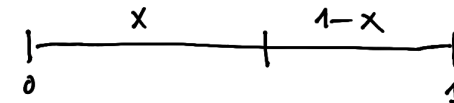


Figure 4.2: Stick breaking visualisation of a beta random variable.

A beta random variable can be visualised as breaking a unit stick of length one into two pieces of length  $x_1 = x$  and  $x_2 = 1 - x$  (Figure 4.2).

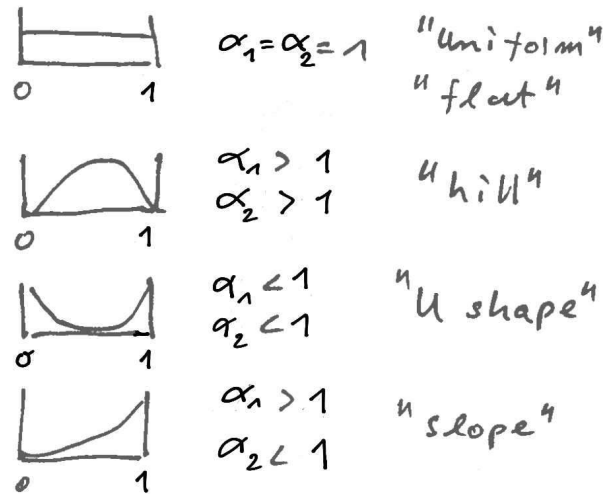


Figure 4.3: Shapes of the pdf of the beta distribution.

Thus, the  $x_i$  may be used as the exclusive proportions or probabilities for  $K = 2$  classes.

The mean is

$$E(x) = E(x_1) = \frac{\alpha_1}{m}$$

and hence

$$E(1 - x) = E(x_2) = \frac{\alpha_2}{m}$$

The variance is

$$\text{Var}(x) = \text{Var}(x_1) = \text{Var}(x_2) = \frac{\alpha_1 \alpha_2}{m^2(m + 1)}$$

The pdf of the beta distribution  $\text{Beta}(\alpha_1, \alpha_2)$  is

$$p(x|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} x^{\alpha_1-1} (1-x)^{\alpha_2-1}$$

This depends on the beta function with arguments  $\alpha_1$  and  $\alpha_2$  defined as

$$B(\alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(m)}$$

The beta distribution can assume a number of different shapes, depending on the values of  $\alpha_1$  and  $\alpha_2$  (see Figure 4.3).

#### R code

The pdf of the beta distribution is given by `dbeta()`, the distribution function is `pbeta()` and the quantile function is `qbeta()`. The corresponding random number generator is `rbeta()`.

### Mean parametrisation

Instead of employing  $\alpha_1$  and  $\alpha_2$  as parameters another useful reparametrisation of the beta distribution is in terms of a mean parameter  $\mu \in [0, 1]$  and a concentration parameter  $m > 0$  so that

$$x \sim \text{Beta}(\alpha_1 = m\mu, \alpha_2 = m(1 - \mu))$$

The concentration and mean parameters can be obtained from  $\alpha_1$  and  $\alpha_2$  by  $m = \alpha_1 + \alpha_2$  and  $\mu = \alpha_1/m$ .

The mean and variance of the beta distribution expressed in terms of  $\mu$  and  $m$  are

$$E(x) = \mu$$

and

$$\text{Var}(x) = \frac{\mu(1 - \mu)}{m + 1}$$

With increasing concentration parameter  $m$  the variance decreases and thus the probability mass becomes more concentrated around the mean.

### Special case: symmetric beta distribution

For  $\alpha_1 = \alpha_2 = \alpha$  the beta distribution becomes the **symmetric beta distribution** with a single shape parameter  $\alpha > 0$ . In mean parametrisation the symmetric beta distribution corresponds to  $\mu = 1/2$  and  $m = 2\alpha$ .

### Special case: uniform distribution

For  $\alpha_1 = \alpha_2 = 1$  the beta distribution becomes the **uniform distribution over the unit interval** with pdf  $p(x) = 1$ . In mean parametrisation the uniform distribution corresponds to  $\mu = 1/2$  and  $m = 2$ .

### 4.3 Normal distribution

The **normal distribution**  $N(\mu, \sigma^2)$  is the most important continuous probability distribution. It is also called **Gaussian distribution** named after [Carl Friedrich Gauss \(1777–1855\)](#).

Special cases are the **standard normal distribution**  $N(0, 1)$  and the **delta distribution**  $\delta$ .

#### Standard parametrisation

The univariate normal distribution  $N(\mu, \sigma^2)$  has two parameters  $\mu$  (location) and  $\sigma^2 > 0$  (variance) and support  $x \in \mathbb{R}$ .

If a random variable  $x$  is normally distributed we write

$$x \sim N(\mu, \sigma^2)$$

with mean

$$E(x) = \mu$$

and variance

$$\text{Var}(x) = \sigma^2$$

The pdf is given by

$$\begin{aligned} p(x|\mu, \sigma^2) &= (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\ &= (\sigma^2)^{-1/2} (2\pi)^{-1/2} e^{-\Delta^2/2} \end{aligned}$$

Here  $\Delta^2 = (x - \mu)^2 / \sigma^2$  is the squared distance between  $x$  and  $\mu$  weighted by the variance  $\sigma^2$ , also known as **squared Mahalanobis distance**.

The normal distribution is sometimes also used by specifying the precision  $1/\sigma^2$  instead of the variance  $\sigma^2$ .

#### R code

The normal pdf is given by `dnorm()`, the distribution function is `pnorm()` and the quantile function is `qnorm()`. The corresponding random number generator is `rnorm()`.

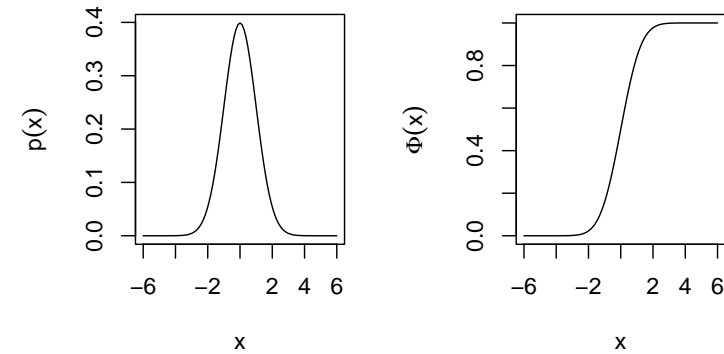


Figure 4.4: Probability density function (left) and cumulative density function (right) of the standard normal distribution.

#### Scale parametrisation

Instead of the variance parameter  $\sigma^2$  it is often also convenient to use the standard deviation  $\sigma = \sqrt{\sigma^2} > 0$  as scale parameter. Similarly, instead of the precision  $1/\sigma^2$  one may wish to use the inverse standard deviation  $w = 1/\sigma$ .

The scale parametrisation is central for location-scale transformations (see below).

#### Special case: standard normal distribution

The **standard normal distribution**  $N(0, 1)$  has mean  $\mu = 0$  and variance  $\sigma^2 = 1$ . The corresponding pdf is

$$p(x) = (2\pi)^{-1/2} e^{-x^2/2}$$

with the squared Mahalanobis distance reduced to  $\Delta^2 = x^2$ .

The cumulative distribution function (cdf) of the standard normal  $N(0, 1)$  is

$$\Phi(x) = \int_{-\infty}^x p(x'|\mu = 0, \sigma^2 = 1) dx'$$

There is no analytic expression for  $\Phi(x)$ . The inverse  $\Phi^{-1}(p)$  is called the quantile function of the standard normal distribution.

Figure 4.4 shows the pdf and cdf of the standard normal distribution.

### Special case: delta distribution

The **delta distribution**  $\delta$  is obtained as the limit of  $N(0, \varepsilon\sigma^2)$  for  $\varepsilon \rightarrow 0$  and where  $\sigma^2$  is a positive number (e.g.  $\sigma^2 = 1$ ). Thus  $\delta$  is a distribution that behaves like an infinite spike at zero.

The corresponding pdf  $\delta(x)$  is called the **Dirac delta function**, even though it is not an ordinary function. It satisfies  $\delta(x) = 0$  for all  $x \neq 0$  and integrates to one, thus representing a point mass at zero.

### Location-scale transformation

Let  $\sigma > 0$  be the positive square root of the variance  $\sigma^2$  and  $w = 1/\sigma$ .

If  $x \sim N(\mu, \sigma^2)$  then  $y = w(x - \mu) \sim N(0, 1)$ . This location-scale transformation corresponds to centring and standardisation of a normal random variable, reducing it to a standard normal random variable.

Conversely, if  $y \sim N(0, 1)$  then  $x = \mu + \sigma y \sim N(\mu, \sigma^2)$ . This location-scale transformation generates the normal distribution from the standard normal distribution.

### Convolution property

The convolution of  $n$  independent, but not necessarily identical, normal distributions results in another normal distribution with corresponding mean and variance:

$$\sum_{i=1}^n N(\mu_i, \sigma_i^2) \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

Hence, any normal random variable can be constructed as the sum of  $n$  suitable independent normal random variables.

Since  $n$  is an arbitrary positive integer the normal distribution is said to be **infinitely divisible**.

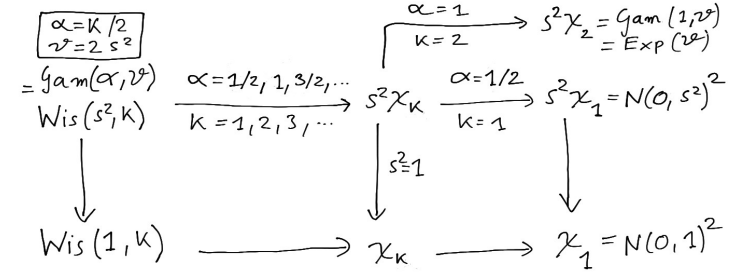


Figure 4.5: The gamma and the univariate Wishart distribution and its relatives.

## 4.4 Gamma distribution

The **gamma distribution**  $\text{Gam}(\alpha, \theta)$  is another widely used continuous distribution and is also known as **univariate Wishart distribution**  $\text{Wis}(s^2, k)$  using a different parametrisation.

It contains as special cases the **scaled chi-squared distribution**  $s^2 \chi_k^2$  (two parameter restrictions) as well as the **univariate standard Wishart distribution**  $\text{Wis}(1, k)$ , the **chi-squared distribution**  $\chi_k^2$  and the **exponential distribution**  $\text{Exp}(\theta)$  (one parameter restrictions). Figure 4.5 illustrates the relationship of the gamma and the univariate Wishart distribution with these related distributions.

### Standard parametrisation

The gamma distribution  $\text{Gam}(\alpha, \theta)$  is a continuous distribution with two parameters  $\alpha > 0$  (shape) and  $\theta > 0$  (scale):

$$x \sim \text{Gam}(\alpha, \theta)$$

and support  $x \in [0, \infty[$  with mean

$$E(x) = \alpha\theta$$

and variance

$$\text{Var}(x) = \alpha\theta^2$$

The gamma distribution is also often used with a rate parameter  $\beta = 1/\theta$ . Therefore one needs to pay attention which parametrisation is used.

The pdf is

$$p(x|\alpha, \theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} e^{-x/\theta}$$

#### 💡 R code

The pdf of the gamma distribution is available in the function `dgamma()`, the distribution function is `pgamma()` and the quantile function is `qgamma()`. The corresponding random number generator is `rgamma()`.

### Wishart parametrisation

The gamma distribution is often used with a different set of parameters  $s^2 = \theta/2 > 0$  (scale) and  $k = 2\alpha > 0$  (shape or concentration). In this form it is known as **univariate or one-dimensional Wishart distribution**

$$x \sim \text{Wis}(s^2, k) = \text{Gam}\left(\alpha = \frac{k}{2}, \theta = 2s^2\right)$$

named after [John Wishart \(1898–1954\)](#).

In the above the scale parameter  $s^2$  is scalar and hence the resulting Wishart distribution is univariate. If instead a matrix-valued scale parameter  $S$  is used this yields the multivariate or  $d$ -dimensional Wishart distribution, see Section 5.4.

In the Wishart parametrisation the mean is

$$E(x) = ks^2$$

and the variance

$$\text{Var}(x) = 2ks^4$$

The pdf in terms of  $s^2$  and  $k$  is

$$p(x|s^2, k) = \frac{1}{\Gamma(k/2)(2s^2)^{k/2}} x^{(k-2)/2} e^{-s^{-2}x/2}$$

### Mean parametrisation

Finally, we also often employ the Wishart resp. gamma distribution in **mean parametrisation**

$$x \sim \text{Wis}\left(s^2 = \frac{\mu}{k}, k\right) = \text{Gam}\left(\alpha = \frac{k}{2}, \theta = \frac{2\mu}{k}\right)$$

with parameters  $\mu = ks^2 > 0$  and  $k > 0$ . In this parametrisation the mean is

$$E(x) = \mu$$

and the variance

$$\text{Var}(x) = \frac{2\mu^2}{k}$$

### Special case: univariate standard Wishart distribution

For  $s^2 = 1$  the univariate Wishart distribution reduces to the **univariate standard Wishart distribution**

$$x \sim \text{Wis}(1, k) = \text{Gam}\left(\alpha = \frac{k}{2}, \theta = 2\right)$$

with mean

$$E(x) = k$$

and variance

$$\text{Var}(x) = 2k$$

The pdf is

$$p(x|k) = \frac{1}{\Gamma(k/2)2^{k/2}} x^{(k-2)/2} e^{-x/2}$$

### Special case: scaled chi-squared distribution

If the shape parameter  $k$  of the Wishart distribution  $\text{Wis}(s^2, k)$  is restricted to the positive integers  $k \in \{1, 2, \dots\}$  the Wishart distribution becomes the **scaled chi-squared distribution**  $s^2 \chi_k^2$  where  $k$  is called the degree of freedom.

This is equivalent to restricting the shape parameter  $\alpha$  of the gamma distribution  $\text{Gam}(\alpha = k/2, \theta = 2s^2)$  to  $\alpha \in \{1/2, 1, 3/2, 2, \dots\}$ .

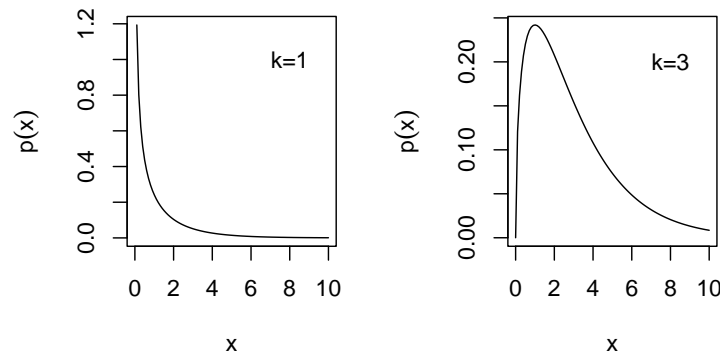


Figure 4.6: Probability density function of the chi-squared distribution.

The scaled chi-squared distribution with  $k = 1$  is the distribution of a squared normal random variable with mean zero. Specifically, if  $z \sim N(0, s^2)$  then  $z^2 \sim s^2 \chi_1^2 = \text{Wis}(s^2, 1) = N(0, s^2)^2$ .

### Special case: chi-squared distribution

If  $k$  is restricted to the positive integers  $k \in \{1, 2, \dots\}$  the univariate standard Wishart distribution reduces to the **chi-squared distribution**

$$x \sim \chi_k^2 = \text{Wis}(s^2 = 1, k) = \text{Gam}\left(\alpha = \frac{k}{2}, \theta = 2\right)$$

where  $k$  is called the degree of freedom.

The chi-squared distribution has mean

$$E(x) = k$$

and variance

$$\text{Var}(x) = 2k$$

Figure 4.6 shows the pdf of the chi-squared distribution for degrees of freedom  $k = 1$  and  $k = 3$ .

The chi-squared distribution with  $k = 1$  is the distribution of a squared standard normal random variable. Specifically, if  $z \sim N(0, 1)$  then  $z^2 \sim \chi_1^2 = \text{Wis}(1, 1) = N(0, 1)^2$ .

#### 💡 R code

The pdf of the chi-squared distribution is given by `dchisq()`. The distribution function is `pchisq()` and the quantile function is `qchisq()`. The corresponding random number generator is `rchisq()`.

### Special case: exponential distribution

If the shape parameter  $\alpha$  of the gamma distribution  $\text{Gam}(\alpha, \theta)$  is set to  $\alpha = 1$ , or if the shape parameter  $k$  of the Wishart distribution  $\text{Wis}(s^2, k)$  is set to  $k = 2$ , we obtain the **exponential distribution**

$$x \sim \text{Exp}(\theta) = \text{Gam}(\alpha = 1, \theta) = \text{Wis}(s^2 = \theta/2, k = 2)$$

with scale parameter  $\theta$ .

It has mean

$$E(x) = \theta$$

and variance

$$\text{Var}(x) = \theta^2$$

and the pdf is

$$p(x|\theta) = \theta^{-1} e^{-x/\theta}$$

Just like the gamma distribution the exponential distribution is also often specified using a rate parameter  $\beta = 1/\theta$  instead of a scale parameter  $\theta$ .

#### 💡 R code

The command `dexp()` returns the pdf of the exponential distribution, `pexp()` is the distribution function and `qexp()` is the quantile function. The corresponding random number generator is `rexp()`.



### Scale transformation

If  $x \sim \text{Gam}(\alpha, \theta)$  then the scaled random variable  $bx$  with  $b > 0$  is also gamma distributed with  $bx \sim \text{Gam}(\alpha, b\theta)$ .

Hence,

- $\theta \text{Gam}(\alpha, 1) = \text{Gam}(\alpha, \theta)$ ,
- $\theta \text{Exp}(1) = \text{Exp}(\theta)$ ,
- $(\mu/k) \text{Wis}(1, k) = \text{Wis}(s^2 = \mu/k, k)$  and
- $s^2 \text{Wis}(1, k) = \text{Wis}(s^2, k)$ .

As  $\chi_k^2$  equals  $\text{Wis}(1, k)$  the last example demonstrates that the **scaled chi-squared distribution**  $s^2 \chi_k^2$  equals the univariate Wishart distribution  $\text{Wis}(s^2, k)$ .

### Convolution property

The convolution of  $n$  gamma distributions with the same scale parameter  $\theta$  but possible different shape parameters  $\alpha_i$  yields another gamma distribution:

$$\sum_{i=1}^n \text{Gam}(\alpha_i, \theta) \sim \text{Gam}\left(\sum_{i=1}^n \alpha_i, \theta\right)$$

Thus, any gamma random variable can be obtained as the sum of  $n$  suitable independent gamma random variables.

In Wishart parametrisation this becomes

$$\sum_{i=1}^n \text{Wis}(s^2, k_i) \sim \text{Wis}\left(s^2, \sum_{i=1}^n k_i\right)$$

As a result, since  $n$  is an arbitrary positive integer, the gamma resp. univariate Wishart distribution is **infinitely divisible**.

The above includes the following two specific constructions:

- If  $x_1, \dots, x_n \sim \text{Exp}(\theta)$  are independent samples from  $\text{Exp}(\theta)$  then the sum  $y = \sum_{i=1}^n x_i \sim \text{Gam}(\alpha = n, \theta)$  is gamma distributed with the same scale parameter.

- The sum of  $k$  independent scaled chi-squared random variables  $s^2 \chi_1^2$  with one degree of freedom and identical scale parameter  $s^2$  yields a scaled chi-squared random variable  $s^2 \chi_k^2$  with degree of freedom  $k$  and the same scale parameter. Thus, if  $z_1, z_2, \dots, z_k \sim N(0, 1)$  are  $k$  independent samples from  $N(0, 1)$  then  $\sum_{i=1}^k z_i^2 \sim \chi_k^2$ .

## 4.5 Inverse gamma distribution

The **inverse gamma distribution**  $\text{IG}(\alpha, \beta)$  is a continuous distribution and is also known as **univariate inverse Wishart distribution**  $\text{IW}(\psi, k)$  using a different parametrisation. It is linked to the gamma distribution  $\text{Gam}(\alpha, \theta)$  aka univariate Wishart distribution  $\text{Wis}(s^2, k)$  (Section 4.4).

Special cases include the **inverse chi-squared distribution**  $\text{Inv-}\chi_k^2$  and the **scaled inverse chi-squared distribution**  $s^2 \text{Inv-}\chi_k^2$ .

### Standard parametrisation

A random variable  $x$  following an **inverse gamma distribution** is denoted by

$$x \sim \text{IG}(\alpha, \beta)$$

with two parameters  $\alpha > 0$  (shape parameter) and  $\beta > 0$  (scale parameter) and support  $x > 0$ .

The mean of the inverse gamma distribution is (for  $\alpha > 1$ )

$$E(x) = \frac{\beta}{\alpha - 1}$$

and the variance (for  $\alpha > 2$ )

$$\text{Var}(x) = \frac{\beta^2}{(\alpha - 1)^2(\alpha - 2)}$$

The inverse gamma distribution  $\text{IG}(\alpha, \beta)$  has pdf

$$p(x|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (1/x)^{\alpha+1} e^{-\beta/x}$$

The inverse gamma distribution and the gamma distribution are directly linked. If  $x \sim \text{IG}(\alpha, \beta)$  then the inverse of  $x$  is gamma distributed with inverted scale parameter

$$\frac{1}{x} \sim \text{Gam}(\alpha, \theta = \beta^{-1})$$

where  $\alpha$  is the shape parameter and  $\theta$  the scale parameter of the gamma distribution.

#### 💡 R code

The `extraDistr` package implements the inverse gamma distribution. The function `extraDistr::dinvgamma()` provides the pdf, `extraDistr::pinvgamma()` the distribution function and `extraDistr::qinvgamma()` is the quantile function. The corresponding random number generator is `extraDistr::rinvgamma()`.

### Wishart parametrisation

The inverse gamma distribution is frequently used with a different set of parameters  $\psi = 2\beta$  (scale parameter) and  $k = 2\alpha$  (shape parameter). In this form it is called **univariate inverse Wishart distribution**

$$x \sim \text{IW}(\psi, k) = \text{IG}\left(\alpha = \frac{k}{2}, \beta = \frac{\psi}{2}\right)$$

In the above the scale parameter  $\psi$  is scalar and hence the resulting inverse Wishart distribution is univariate. If instead a matrix-valued scale parameter  $\Psi$  is used this yields the multivariate or  $d$ -dimensional inverse Wishart distribution, see Section 5.5.

In the Wishart parametrisation the mean is (for  $k > 2$ )

$$\text{E}(x) = \frac{\psi}{k-2}$$

and the variance is (for  $k > 4$ )

$$\text{Var}(x) = \frac{2\psi^2}{(k-4)(k-2)^2}$$

The pdf in terms of  $\psi$  and  $k$  is

$$p(x|\psi, k) = \frac{(\psi/2)^{(k/2)}}{\Gamma(k/2)} x^{-(k+2)/2} e^{-\psi x^{-1}/2}$$

The univariate inverse Wishart and the univariate Wishart distributions are linked. If  $x \sim \text{IW}(\psi, k)$  then the inverse of  $x$  is Wishart distributed with inverted scale parameter:

$$\frac{1}{x} \sim \text{Wis}(s^2 = \psi^{-1}, k)$$

where  $k$  is shape parameter and  $s^2$  the scale parameter of the Wishart distribution.

### Mean parametrisation

Instead of  $\psi$  and  $k$  we may also equivalently use  $\mu = \psi/(\nu - 2)$  and  $\kappa = \nu - 2$  as parameters for the univariate inverse Wishart distribution, so that

$$x \sim \text{IW}(\psi = \kappa\mu, k = \kappa + 2) = \text{IG}\left(\alpha = \frac{\kappa + 2}{2}, \beta = \frac{\mu\kappa}{2}\right)$$

has mean (for  $\kappa > 0$ )

$$\text{E}(x) = \mu$$

and the variance (for  $\kappa > 2$ )

$$\text{Var}(x) = \frac{2\mu^2}{\kappa - 2}$$

The **mean parametrisation** is useful in Bayesian analysis when employing the inverse gamma aka univariate inverse Wishart distribution as prior and posterior distribution.

### Biased mean parametrisation

Using  $\tau^2 = \frac{\psi}{k}$  as biased mean parameter together with  $\nu = k$  we arrive at the **biased mean parametrisation**

$$x \sim \text{IW}(\psi = \nu\tau^2, k = \nu) = \text{IG}\left(\alpha = \frac{\nu}{2}, \beta = \frac{\nu\tau^2}{2}\right)$$

with mean (for  $\nu > 2$ )

$$E(x) = \frac{\nu}{\nu - 2} \tau^2 = \mu$$

and variance ( $\nu > 4$ )

$$\text{Var}(x) = \left( \frac{\nu}{\nu - 2} \right)^2 \frac{2\tau^4}{\nu - 4}$$

As  $\tau^2 = \mu(\nu - 2)/\nu$  for large  $\nu$  the parameter  $\tau^2$  will become identical to the true mean  $\mu$ .

This parametrisation is useful to derive the location-scale  $t$ -distribution with its matching parameters (see Section 4.6). It is also common in Bayesian analysis.

### Special case: inverse chi-squared distribution

If the scale parameter in  $\text{IW}(\psi, k)$  is set to  $\psi = 1$  and  $k$  is restricted to the positive integers  $\{1, 2, \dots\}$  then the univariate inverse Wishart distribution reduces to the **inverse chi-squared distribution**

$$x \sim \text{Inv-}\chi_k^2 = \text{IW}(\psi = 1, k) = \text{IG}\left(\alpha = \frac{k}{2}, \beta = \frac{1}{2}\right)$$

where  $k$  is called the degree of freedom.

The inverse chi-squared distribution has mean (for  $k > 2$ )

$$E(x) = \frac{1}{k - 2}$$

and the variance is (for  $k > 4$ )

$$\text{Var}(x) = \frac{2}{(k - 2)^2(k - 4)}$$

The inverse chi-squared distribution and the chi-squared distribution are linked. If  $x \sim \text{Inv-}\chi_k^2$  then the inverse of  $x$  is chi-squared distributed:

$$\frac{1}{x} \sim \chi_k^2$$

where  $k$  is the degree of freedom.

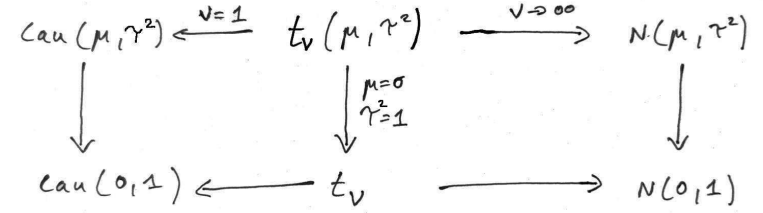


Figure 4.7: The location-scale  $t$ -distribution and its relatives.

### Scale transformation

If  $x \sim \text{IG}(\alpha, \beta)$  then the scaled random variable  $bx$  with  $b > 0$  is also inverse gamma distributed with  $bx \sim \text{IG}(\alpha, b\beta)$ .

Hence,

- $\beta \text{IG}(\alpha, 1) = \text{IG}(\alpha, \beta)$ ,
- $\psi \text{IW}(1, k) = \text{IW}(\psi, k)$ ,
- $\kappa \mu \text{IW}(1, k = \kappa + 2) = \text{IW}(\psi = \kappa \mu, k = \kappa + 2)$  and
- $\nu \tau^2 \text{IW}(1, k = \nu) = \text{IW}(\psi = \nu \tau^2, k = \nu)$

As  $\text{Inv-}\chi_k^2$  equals  $\text{IW}(\psi = 1, k)$  the **scaled inverse chi-squared distribution**  $\psi \text{Inv-}\chi_k^2$  is thus equivalent to the univariate inverse Wishart distribution  $\text{IW}(\psi, k)$ . If a random variable follows the scaled inverse chi-squared distribution its inverse follows the corresponding scaled chi-squared distribution. Specifically, if  $x \sim \psi \text{Inv-}\chi_k^2$  then  $1/x \sim \psi^{-1} \chi_k^2$ .

The scaled inverse chi-squared distribution is frequently used in the biased mean parametrisation with  $\tau = \psi/\nu$  and  $\nu = k$ . Then  $\psi \text{Inv-}\chi_k^2$  is equal to  $\nu \tau^2 \text{Inv-}\chi_\nu^2 = \text{IW}(\psi = \nu \tau^2, k = \nu)$  which is sometimes also written as  $\text{Inv-}\chi_\nu^2(\nu, \tau^2)$ . If  $x \sim \nu \tau^2 \text{Inv-}\chi_\nu^2$  then  $1/x \sim 1/(\nu \tau^2) \chi_\nu^2$ .

## 4.6 Location-scale $t$ -distribution

The **location-scale  $t$ -distribution**  $t_\nu(\mu, \tau^2)$  is a continuous distribution and is a generalisation of the normal distribution  $N(\mu, \tau^2)$  (Section 4.3) with an additional parameter  $\nu > 0$  (degrees of freedom) controlling the probability mass in the tails.

Special cases include the **Student's  $t$ -distribution**  $t_\nu$ , the **normal distribution**  $N(\mu, \tau^2)$  and the **Cauchy distribution**  $\text{Cau}(\mu, \tau^2)$ . Figure 4.7

illustrates the relationship of the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  with these related distributions.

### Standard parametrisation

If a random variable  $x \in \mathbb{R}$  follows the **location-scale  $t$ -distribution** we write

$$x \sim t_\nu(\mu, \tau^2)$$

where  $\mu$  is the location and  $\tau^2$  the dispersion parameter. The parameter  $\nu > 0$  prescribes the degrees of freedom. For small values of  $\nu$  the distribution is heavy-tailed and as a result only moments of order smaller than  $\nu$  are finite and defined.

The mean is (for  $\nu > 1$ )

$$E(x) = \mu$$

and the variance (for  $\nu > 2$ )

$$\text{Var}(x) = \frac{\nu}{\nu - 2} \tau^2$$

The pdf of  $t_\nu(\mu, \tau^2)$  is

$$p(x|\mu, \tau^2, \nu) = (\tau^2)^{-1/2} \frac{\Gamma(\frac{\nu+1}{2})}{(\pi\nu)^{1/2} \Gamma(\frac{\nu}{2})} \left(1 + \frac{\Delta^2}{\nu}\right)^{-(\nu+1)/2}$$

with  $\Delta^2 = (x - \mu)^2 / \tau^2$  the squared Mahalanobis distance between  $x$  and  $\mu$ .

#### R code

The package `extraDistr` implements the location-scale  $t$ -distribution. The function `extraDistr::dlst()` returns the pdf, `extraDistr::plst()` the distribution function and `extraDistr::qlst()` is the quantile function. The corresponding random number generator is `extraDistr::rlst()`.

### Scale parametrisation

Instead of the dispersion parameter  $\tau^2$  it is often also convenient to use the scale parameter  $\tau = \sqrt{\tau^2} > 0$ . Similarly, instead of the inverse dispersion  $1/\tau^2$  one may wish to use the inverse scale  $w = 1/\tau$ .

The scale parametrisation is central for location-scale transformations (see below).

### Special case: Student's $t$ -distribution

With  $\mu = 0$  and  $\tau^2 = 1$  the location-scale  $t$ -distribution reduces to the **standard  $t$ -distribution**  $t_\nu = t_\nu(0, 1)$ . It is commonly known **Student's  $t$ -distribution** named after “Student” which was the pseudonym of [William Sealy Gosset \(1876–1937\)](#). It is a generalisation of the standard normal distribution  $N(0, 1)$  to allow for heavy tails.

The distribution has mean  $E(x) = 0$  (for  $\nu > 1$ ) and variance  $\text{Var}(x) = \frac{\nu}{\nu-2}$  (for  $\nu > 2$ ).

The pdf of  $t_\nu$  is

$$p(x|\nu) = \frac{\Gamma(\frac{\nu+1}{2})}{(\pi\nu)^{1/2} \Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-(\nu+1)/2}$$

with the squared Mahalanobis distance reducing to  $\Delta^2 = x^2$ .

#### R code

The command `dt()` returns the pdf of the  $t$ -distribution, `pt()` is distribution function and `qt()` the quantile function. The corresponding random number generator is `rt()`.

### Special case: normal distribution

For  $\nu \rightarrow \infty$  the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  reduces to the **normal distribution**  $N(\mu, \tau^2)$  (Section 4.3). Correspondingly, for  $\nu \rightarrow \infty$  the Student's  $t$ -distribution becomes equal to the **standard normal distribution**  $N(0, 1)$ .

See Section 5.6 for further details.

### Special case: Cauchy distribution

For  $\nu = 1$  the location-scale  $t$ -distribution becomes the **Cauchy distribution**  $\text{Cau}(\mu, \tau^2) = t_1(\mu, \tau^2)$  named after [Augustin-Louis Cauchy \(1789–1857\)](#).

Its mean, variance and other higher moments are all undefined.

It has pdf

$$p(x|\mu, \tau^2) = (\tau^2)^{-1/2}(\pi(1 + \Delta^2))^{-1} \\ = \frac{\tau}{\pi(\tau^2 + (x - \mu)^2)}$$

with  $\tau = \sqrt{\tau^2} > 0$ .

Note that in the above we employ  $\tau^2$  as dispersion parameter as this parallels the location-scale  $t$ -distribution and the normal distribution but very often the Cauchy distribution is used with  $\tau > 0$  as scale parameter.

#### 💡 R code

The command `dcauchy()` returns the pdf of the Cauchy distribution, `pcachy()` is the distribution function and `qcauchy()` the quantile function. The corresponding random number generator is `rcachy()`.

### Special case: standard Cauchy distribution

The **standard Cauchy distribution**  $\text{Cau}(0, 1) = t_1(0, 1) = t_1$  is obtained by setting  $\mu = 0$  and  $\tau^2 = 1$  (Cauchy distribution) or, equivalently, by setting  $\nu = 1$  (Student's  $t$ -distribution).

It has pdf

$$p(x) = \frac{1}{\pi(1 + x^2)}$$

### Location-scale transformation

Let  $\tau > 0$  be the positive square root of  $\tau^2$  and  $w = 1/\tau$ .

If  $x \sim t_\nu(\mu, \tau^2)$  then  $y = w(x - \mu) \sim t_\nu$ . This location-scale transformation reduces a location-scale  $t$ -distributed random variable to a Student's  $t$ -distributed random variable.

Conversely, if  $y \sim t_\nu$  then  $x = \mu + \tau y \sim t_\nu(\mu, \tau^2)$ . This location-scale transformation generates the location-scale  $t$ -distribution from the Student's  $t$ -distribution.

For the special case of the Cauchy distribution (corresponding to  $\nu = 1$ ) similar relations hold between it and the standard Cauchy distribution. If  $x \sim \text{Cau}(\mu, \tau^2)$  then  $y = w(x - \mu) \sim \text{Cau}(0, 1)$ . Conversely, if  $y \sim \text{Cau}(0, 1)$  then  $x = \mu + \tau y \sim \text{Cau}(\mu, \tau^2)$ .

### Convolution property

The location-scale  $t$ -distribution is not generally closed under convolution, with the exception of two special cases, the normal distribution ( $\nu = \infty$ ), see Section 4.3, and the Cauchy distribution ( $\nu = 1$ ).

For the Cauchy distribution with  $\tau_i^2 = a_i^2 \tau^2$ , where  $a_i > 0$  are positive scalars,

$$\sum_{i=1}^n \text{Cau}(\mu_i, a_i^2 \tau^2) \sim \text{Cau}\left(\sum_{i=1}^n \mu_i, \left(\sum_{i=1}^n a_i\right)^2 \tau^2\right)$$

### Location-scale $t$ -distribution as compound distribution

The location-scale  $t$ -distribution can be obtained as mixture of normal distributions with identical mean and varying variance. Specifically, let  $z$  be a univariate inverse Wishart random variable

$$z \sim \text{IW}(\psi = \nu, k = \nu) = \text{IG}\left(\alpha = \frac{\nu}{2}, \beta = \frac{\nu}{2}\right)$$

and let  $x|z$  be normal

$$x|z \sim N(\mu, \sigma^2 = z\tau^2)$$

Then the resulting marginal (scale mixture) distribution for  $x$  is the location-scale  $t$ -distribution

$$x \sim t_\nu(\mu, \tau^2)$$

An alternative way to arrive at  $t_\nu(\mu, \tau^2)$  is to include  $\tau^2$  as parameter in the inverse Wishart distribution

$$z \sim \tau^2 \text{IW}(\psi = \nu, k = \nu) = \text{IW}(\psi = \nu\tau^2, k = \nu)$$

and let

$$x|z \sim N(\mu, \sigma^2 = z)$$

Note that  $\tau^2$  is now the biased mean parameter of the univariate inverse Wishart distribution. This characterisation is useful in Bayesian analysis.

## 5 Multivariate distributions

### 5.1 Multinomial distribution

The **multinomial distribution**  $\text{Mult}(n, \theta)$  is the multivariate generalisation of the binomial distribution  $\text{Bin}(n, \theta)$  (Section 4.1) from two classes to  $K$  classes.

A special case is the **categorical distribution**  $\text{Cat}(\theta)$  that generalises the Bernoulli distribution  $\text{Ber}(\theta)$ .

#### Standard parametrisation

A multinomial random variable  $x$  describes the allocation of  $n$  items to  $K$  classes. We write

$$x \sim \text{Mult}(n, \theta)$$

where the parameter vector  $\theta = (\theta_1, \dots, \theta_K)^T$  specifies the probability of each of  $K$  classes, with  $\theta_i \in [0, 1]$  and  $\theta^T \mathbf{1}_K = \sum_{i=1}^K \theta_i = 1$ . Thus there are  $K - 1$  independent elements in  $\theta$ . The number of classes  $K$  is implicitly given by the dimension of the vector  $\theta$ . Each element of the vector  $x = (x_1, \dots, x_K)^T$  is an integer  $x_i \in \{0, 1, \dots, n\}$  and  $x$  satisfies the constraint  $x^T \mathbf{1}_K = \sum_{i=1}^K x_i = n$ . Therefore the support of  $x$  is a  $K - 1$  dimensional space and it notably depends on  $n$ .

The multinomial distribution is best illustrated by an urn model distributing  $n$  items into  $K$  bins where  $\theta$  contains the corresponding bin probabilities (Figure 5.1).

The expected value is

$$E(x) = n\theta$$

The covariance matrix is

$$\text{Var}(x) = n(\text{Diag}(\theta) - \theta\theta^T)$$

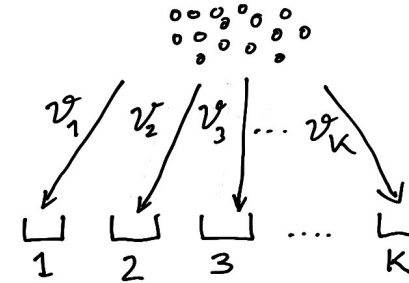


Figure 5.1: Multinomial urn model.

The covariance matrix is singular by construction because of the dependencies among the elements of  $x$ .

The corresponding pmf is

$$p(x|\theta) = \binom{n}{x_1, \dots, x_K} \prod_{i=1}^K \theta_i^{x_i}$$

where  $\binom{n}{x_1, \dots, x_K}$  is the multinomial coefficient. While all  $K$  elements of  $x$  appear in the pmf recall that due to the dependencies among the  $x_i$  the pmf is defined over a  $K - 1$  dimensional support.

For  $K = 2$  the multinomial distribution reduces to the binomial distribution (Section 4.1).

#### 💡 R code

The pmf of the multinomial distribution is given by `dmultinom()`.  
The corresponding random number generator is `rmultinom()`.

#### Mean parametrisation

Instead of  $\theta$  one may also use a mean parameter  $\mu$ , with elements  $\mu_i \in [0, n]$  and  $\mu^T \mathbf{1}_K = \sum_{i=1}^K \mu_i = n$ , so that

$$x \sim \text{Mult}\left(n, \theta = \frac{\mu}{n}\right)$$

The mean parameter  $\mu$  can be obtained from  $\theta$  and  $n$  by  $\mu = n\theta$ . Note that the parameter space for  $\mu$  and the support of  $x$  are both of dimension  $K - 1$ .

The mean and variance of the multinomial distribution expressed in terms of  $\mu$  and  $n$  are

$$E(x) = \mu$$

and

$$\text{Var}(x) = \text{Diag}(\mu) - \frac{\mu\mu^T}{n}$$

The covariance matrix is singular by construction because of the dependencies among the elements of  $x$ .

### Special case: categorical distribution

For  $n = 1$  the multinomial distribution reduces to the **categorical distribution**  $\text{Cat}(\theta)$  which in turn is the multivariate generalisation of the Bernoulli distribution  $\text{Ber}(\theta)$  from two classes to  $K$  classes.

If a random variable  $x$  follows the categorical distribution we write

$$x \sim \text{Cat}(\theta)$$

with class probabilities  $\theta$  and  $\theta^T \mathbf{1}_K = 1$ . The support is  $x_i \in \{0, 1\}$  and  $x^T \mathbf{1}_K = 1$  and is a  $K - 1$  dimensional space.

The random vector  $x$  takes the form of an indicator vector  $x = (x_1, \dots, x_K)^T = (0, 0, \dots, 1, \dots, 0)^T$  containing zeros everywhere except for a single element  $x_k = 1$  indicating the class  $k$  to which the item has been allocated. This is called “one hot encoding”, as opposed to “integer encoding”, i.e. stating the class number  $k$ .

The expected value is

$$E(x) = \theta$$

The covariance matrix is

$$\text{Var}(x) = \text{Diag}(\theta) - \theta\theta^T$$

The covariance matrix is singular by construction because of the dependencies among the elements of  $x$ . This follows directly from the definition of the variance  $\text{Var}(x) = E(xx^T) - E(x)E(x)^T$  and noting that  $x_i^2 = x_i$  and  $x_i x_j = 0$  if  $i \neq j$ .

The corresponding pmf is

$$p(x|\theta) = \prod_{k=1}^K \theta_k^{x_k} = \begin{cases} \theta_k & \text{if } x_k = 1 \end{cases}$$

Recall that the pmf is defined over the  $K - 1$  dimensional support of  $x$ .

For  $K = 2$  the categorical distribution reduces to the Bernoulli  $\text{Ber}(\theta)$  distribution, with  $\theta_1 = \theta$ ,  $\theta_2 = 1 - \theta$  and  $x_1 = x$  and  $x_2 = 1 - x$ .

### Convolution property

The convolution of  $n$  multinomial distributions, each with identical bin probabilities  $\theta$  but possibly different number of items  $n_i$ , yields another multinomial distribution with the same parameter  $\theta$ :

$$\sum_{i=1}^n \text{Mult}(n_i, \theta) \sim \text{Mult}\left(\sum_{i=1}^n n_i, \theta\right)$$

It follows that the multinomial distribution with  $n$  items is the result of the convolution of  $n$  categorical distributions:

$$\sum_{i=1}^n \text{Cat}(\theta) \sim \text{Mult}(n, \theta)$$

Thus, repeating the same categorical trial  $n$  times and counting the total number of allocations in each bin yields a multinomial random variable.

## 5.2 Dirichlet distribution

The **Dirichlet distribution**  $\text{Dir}(\alpha)$  is the multivariate generalisation of the beta distribution  $\text{Beta}(\alpha_1, \alpha_2)$  (Section 4.2) that is useful to model proportions or probabilities for  $K \geq 2$  classes. It is named after [Peter Gustav Lejeune Dirichlet \(1805–1859\)](#).

It includes the **uniform distribution** over the  $K - 1$  unit simplex as special case.

### Standard parametrisation

A Dirichlet distributed random vector is denoted by

$$\mathbf{x} \sim \text{Dir}(\boldsymbol{\alpha})$$

with shape parameter  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)^T > 0$  and  $K \geq 2$ . Let  $m = \boldsymbol{\alpha}^T \mathbf{1}_K = \sum_{i=1}^K \alpha_i$ . The support of  $\mathbf{x}$  is the  $K - 1$  dimensional unit simplex given by  $x_i \in [0, 1]$  and  $\mathbf{x}^T \mathbf{1}_K = \sum_{i=1}^K x_i = 1$ . Thus, the Dirichlet distribution is defined over a  $K - 1$  dimensional space.

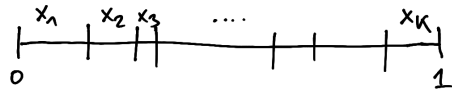


Figure 5.2: Stick breaking visualisation of a Dirichlet random variable.

A Dirichlet random variable can be visualised as breaking a unit stick into  $K$  individual pieces of lengths  $x_1, x_2, \dots, x_K$  adding up to one (Figure 5.2). Thus, the  $x_i$  may be used as the exclusive proportions or probabilities for  $K$  classes.

The mean is

$$E(\mathbf{x}) = \frac{\boldsymbol{\alpha}}{m}$$

and the variance is

$$\text{Var}(\mathbf{x}) = \frac{m \text{Diag}(\boldsymbol{\alpha}) - \boldsymbol{\alpha} \boldsymbol{\alpha}^T}{m^2(m+1)}$$

The covariance matrix is singular by construction because of the dependencies among the elements of  $\mathbf{x}$ . In component notation it is

$$\text{Cov}(x_i, x_j) = \frac{[i=j]m\alpha_i - \alpha_i\alpha_j}{m^2(m+1)}$$

where the indicator function  $[i=j]$  equals 1 if  $i=j$  and 0 otherwise.

The pdf of the Dirichlet distribution  $\text{Dir}(\boldsymbol{\alpha})$  is

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K x_k^{\alpha_k-1}$$

This depends on the beta function with multivariate argument  $\boldsymbol{\alpha}$  defined as

$$B(\boldsymbol{\alpha}) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(m)}$$

While all  $K$  elements of  $\mathbf{x}$  appear in the pdf recall that due the dependencies among the  $x_i$  the pdf is defined over a  $K - 1$  dimensional support.

For  $K = 2$  the Dirichlet distribution reduces to the beta distribution (Section 4.2).

#### 💡 R code

The `extraDistr` package implements the Dirichlet distribution. The pmf of the Dirichlet distribution is given by `extraDistr::ddirichlet()`. The corresponding random number generator is `extraDistr::rdirichlet()`.

### Mean parametrisation

Instead of employing  $\boldsymbol{\alpha}$  as parameter vector another useful reparametrisation of the Dirichlet distribution is in terms of a mean parameter  $\boldsymbol{\mu}$ , with elements  $\mu_i \in [0, 1]$  and  $\boldsymbol{\mu}^T \mathbf{1}_K = \sum_{i=1}^K \mu_i = 1$ , and a concentration parameter  $m > 0$  so that

$$\mathbf{x} \sim \text{Dir}(\boldsymbol{\alpha} = m\boldsymbol{\mu})$$

The concentration and mean parameters can be obtained from  $\boldsymbol{\alpha}$  by  $m = \boldsymbol{\alpha}^T \mathbf{1}_K$  and  $\boldsymbol{\mu} = \boldsymbol{\alpha}/m$ . The space of possible values for the mean parameter  $\boldsymbol{\mu}$  and the support of  $\mathbf{x}$  are both of dimension  $K - 1$ .

The mean and variance of the Dirichlet distribution expressed in terms of  $\boldsymbol{\mu}$  and  $m$  are

$$E(\mathbf{x}) = \boldsymbol{\mu}$$

and

$$\text{Var}(\mathbf{x}) = \frac{\text{Diag}(\boldsymbol{\mu}) - \boldsymbol{\mu} \boldsymbol{\mu}^T}{m+1}$$



The covariance matrix is singular by construction because of the dependencies among the elements of  $x$ . In component notation it is

$$\begin{aligned}\text{Cov}(x_i, x_j) &= \frac{[i=j]\mu_i - \mu_i\mu_j}{m+1} \\ &= \begin{cases} \mu_i(1-\mu_i)/(m+1) & \text{if } i=j \\ -\mu_i\mu_j/(m+1) & \text{if } i \neq j \end{cases}\end{aligned}$$

### Special case: symmetric Dirichlet distribution

For  $\alpha = \alpha \mathbf{1}_K$  the Dirichlet distribution becomes the **symmetric beta distribution** with a single shape parameters  $\alpha > 0$ . In mean parametrisation the symmetric Dirichlet distribution corresponds to  $\mu = \mathbf{1}_K/K$  and  $m = \alpha K$ .

### Special case: uniform distribution

For  $\alpha = \mathbf{1}_K$  the Dirichlet distribution becomes the uniform distribution over the  $K-1$  unit simplex with pdf  $p(x) = 1/\Gamma(K)$ . In mean parametrisation the symmetric Dirichlet distribution corresponds to  $\mu = \mathbf{1}_K/K$  and  $m = K$ .

## 5.3 Multivariate normal distribution

The **multivariate normal distribution**  $N(\mu, \Sigma)$  generalises the univariate normal distribution  $N(\mu, \sigma^2)$  (Section 4.3) from one to  $d$  dimensions.

Special cases are the **multivariate standard normal distribution**  $N(\mathbf{0}, I)$  and the **multivariate delta distribution**  $\delta$ .

### Standard parametrisation

The multivariate normal distribution  $N(\mu, \Sigma)$  has a mean or location parameter  $\mu$  (a  $d$  dimensional vector), a variance parameter  $\Sigma$  (a  $d \times d$  positive definite symmetric matrix) and support  $x \in \mathbb{R}^d$ .

If a random vector  $x = (x_1, x_2, \dots, x_d)^T$  follows a multivariate normal distribution we write

$$x \sim N(\mu, \Sigma)$$

with mean

$$E(x) = \mu$$

and variance

$$\text{Var}(x) = \Sigma$$

In the above notation the dimension  $d$  is implicitly given by the dimensions of  $\mu$  and  $\Sigma$  but for clarity one often also writes  $N_d(\mu, \Sigma)$  to explicitly indicate the dimension.

The pdf is given by

$$\begin{aligned}p(x|\mu, \Sigma) &= \det(2\pi\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right) \\ &= \det(\Sigma)^{-1/2} (2\pi)^{-d/2} e^{-\Delta^2/2}\end{aligned}$$

Here  $\Delta^2 = (x-\mu)^T \Sigma^{-1}(x-\mu)$  is the **squared Mahalanobis distance** between  $x$  and  $\mu$  taking into account the variance  $\Sigma$ . Note that this pdf is a joint pdf over the  $d$  elements  $x_1, \dots, x_d$  of the random vector  $x$ .

The multivariate normal distribution is sometimes also used by specifying the precision matrix  $\Sigma^{-1}$  instead of the variance  $\Sigma$ .

For  $d = 1$  the random vector  $x = x$  is a scalar,  $\mu = \mu$ ,  $\Sigma = \sigma^2$  and thus the multivariate normal distribution reduces to the univariate normal distribution (Section 4.3).

#### R code

The `mnormt` package implements the multivariate normal distribution. The function `mnormt::dmnorm()` provides the pdf and `mnormt::pmnorm()` returns the distribution function. The function `mnormt::rmnorm()` is the corresponding random number generator.

The `mniw` package also implements the multivariate normal distribution. The pdf of the Wishart distribution is given by `mniw::dmNorm()`. The corresponding random number generator is `mniw::rmNorm()`.

### Scale parametrisation

In the univariate case it is straightforward to use the standard deviation  $\sigma$  as scale parameter instead of the variance  $\sigma^2$ , and similarly the inverse standard deviation  $w = 1/\sigma$  instead of the precision  $\sigma^{-2}$ . However, in

the multivariate setting with a matrix variance parameter  $\Sigma$  it is less obvious how to define a suitable matrix scale parameter.

Let  $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$  be the eigendecomposition of the positive definite matrix  $\Sigma$ . Then  $\Sigma^{1/2} = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{U}^T$  is the principal matrix square root and  $\Sigma^{-1/2} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T$  the inverse principal matrix square root. Furthermore, let  $\mathbf{Q}$  be an arbitrary orthogonal matrix with  $\mathbf{Q}^T\mathbf{Q} = \mathbf{Q}\mathbf{Q}^T = \mathbf{I}$ .

Then  $\mathbf{W} = \mathbf{Q}\Sigma^{-1/2}$  is called a **whitening matrix** based on  $\Sigma$  and  $\mathbf{L} = \mathbf{W}^{-1} = \Sigma^{1/2}\mathbf{Q}^T$  is the corresponding **inverse whitening matrix**. By construction, the matrix  $\mathbf{L}$  provides a factorisation of the covariance matrix by  $\mathbf{L}\mathbf{L}^T = \Sigma$ . Similarly,  $\mathbf{W}$  factorises the precision matrix by  $\mathbf{W}^T\mathbf{W} = \Sigma^{-1}$ . The two matrices thus provide the basis for the scale parametrisation of the multivariate normal distribution.

Specifically, the matrix  $\mathbf{L}$  is used in place of  $\Sigma$  and plays the role of the matrix scale parameter (corresponding to  $\sigma$  in the univariate setting) and  $\mathbf{W}$  is used in place of the precision matrix  $\Sigma^{-1}$  and plays the role of the inverse matrix scale parameter (corresponding to  $1/\sigma$  in the univariate case). The determinants occurring in the multivariate normal pdf can be rewritten in terms of  $\mathbf{L}$  and  $\mathbf{W}$  using the identities  $|\det(\mathbf{W})| = \det(\Sigma)^{-1/2}$  and  $|\det(\mathbf{L})| = \det(\Sigma)^{1/2}$  as  $\det(\mathbf{Q}) = \pm 1$ .

Since  $\mathbf{Q}$  can be freely chosen the matrices  $\mathbf{W}$  and  $\mathbf{L}$  are not fully determined by  $\Sigma$  alone but there is rotational freedom due to  $\mathbf{Q}$ . Standard choices are

- $\mathbf{Q}^{\text{ZCA}} = \mathbf{I}$  for ZCA-type factorisation with  $\mathbf{W}^{\text{ZCA}} = \Sigma^{-1/2}$  and
- $\mathbf{Q}^{\text{PCA}} = \mathbf{U}^T$  for PCA-type factorisation with  $\mathbf{W}^{\text{PCA}} = \mathbf{\Lambda}^{-1/2}\mathbf{U}^T$ . Note that the matrix  $\mathbf{U}$  is not unique because its columns (eigenvectors) can have different signs (directions), hence  $\mathbf{W}^{\text{PCA}}$  and  $\mathbf{L}^{\text{PCA}}$  are also not unique without further constraints, such as positive diagonal elements of the (inverse) whitening matrix.
- A third common choice is to compute  $\mathbf{L}$  directly by Cholesky decomposition of  $\Sigma$ , which yields an  $\mathbf{L}^{\text{Chol}}$  (and also a  $\mathbf{W}^{\text{Chol}}$ ) in the form of a lower-triangular matrix with a positive diagonal, and a corresponding underlying  $\mathbf{Q}^{\text{Chol}} = (\mathbf{L}^{\text{Chol}})^T \Sigma^{-1/2}$ .

Finally, the whitening matrix  $\mathbf{W}$  and its inverse may also be constructed from the correlation matrix  $\mathbf{P}$  and the diagonal matrix containing the variances  $\mathbf{V}$  (with  $\Sigma = \mathbf{V}^{1/2}\mathbf{P}\mathbf{V}^{1/2}$ ) in the form  $\mathbf{W} = \mathbf{Q}\mathbf{P}^{-1/2}\mathbf{V}^{-1/2}$  and  $\mathbf{L} = \mathbf{V}^{1/2}\mathbf{P}^{1/2}\mathbf{Q}^T$ .

### Special case: multivariate standard normal distribution

The **multivariate standard normal distribution**  $N(\mathbf{0}, \mathbf{I})$  has mean  $\mu = \mathbf{0}$  and variance  $\Sigma = \mathbf{I}$ . The corresponding pdf is

$$p(\mathbf{x}) = (2\pi)^{-d/2} e^{-\mathbf{x}^T \mathbf{x} / 2}$$

with the squared Mahalanobis distance reduced to  $\Delta^2 = \mathbf{x}^T \mathbf{x} = \sum_{i=1}^d x_i^2$ .

The density of the multivariate standard normal distribution is the product of the corresponding univariate standard normal densities

$$p(\mathbf{x}) = \prod_{i=1}^d (2\pi)^{-1/2} e^{-x_i^2/2}$$

and therefore the elements  $x_i$  of  $\mathbf{x} = (x_1, \dots, x_d)^T$  are independent of each other.

### Special case: multivariate delta distribution

The **multivariate delta distribution**  $\delta$  is obtained as the limit of  $N(\mathbf{0}, \varepsilon \mathbf{A})$  for  $\varepsilon \rightarrow 0$  and where  $\mathbf{A}$  is a positive definite matrix (e.g.  $\mathbf{A} = \mathbf{I}$ ). Thus  $\delta$  is a distribution that behaves like an infinite spike at zero.

The corresponding pdf  $\delta(\mathbf{x})$  is called the **multivariate Dirac delta function**, even though it is not an ordinary function. It satisfies  $\delta(\mathbf{x}) = 0$  for all  $\mathbf{x} \neq \mathbf{0}$  and integrates to one, thus representing a point mass at zero.

### Location-scale transformation

Let  $\mathbf{W}$  be a whitening matrix for  $\Sigma$  and  $\mathbf{L}$  the corresponding inverse whitening matrix.

If  $\mathbf{x} \sim N(\mu, \Sigma)$  then  $\mathbf{y} = \mathbf{W}(\mathbf{x} - \mu) \sim N(\mathbf{0}, \mathbf{I})$ . This location-scale transformation corresponds to centring and whitening (i.e. standardisation and decorrelation) of a multivariate normal random variable.

Conversely, if  $\mathbf{y} \sim N(\mathbf{0}, \mathbf{I})$  then  $\mathbf{x} = \mu + \mathbf{L}\mathbf{y} \sim N(\mu, \Sigma)$ . This location-scale transformation generates the multivariate normal distribution from the multivariate standard normal distribution.

Note that under the location-scale transformation  $\mathbf{x} = \mu + \mathbf{L}\mathbf{y}$  with  $\text{Var}(\mathbf{y}) = \mathbf{I}$  we get  $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{L}$ . This provides a means to choose

between different (inverse) whitening transformation and the corresponding factorisations of  $\Sigma$  and  $\Sigma^{-1}$ . For example, if positive correlation between corresponding elements in  $\mathbf{x}$  and  $\mathbf{y}$  is desired then the diagonal elements in  $\mathbf{L}$  must be positive.

### Convolution property

The convolution of  $n$  independent, but not necessarily identical, multivariate normal distributions of the same dimension  $d$  results in another  $d$ -dimensional multivariate normal distribution with corresponding mean and variance:

$$\sum_{i=1}^n N(\mu_i, \Sigma_i) \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \Sigma_i\right)$$

Hence, any multivariate normal random variable can be constructed as the sum of  $n$  suitable independent multivariate normal random variables.

Since  $n$  is an arbitrary positive integer the multivariate normal distribution is said to be **infinitely divisible**.

## 5.4 Wishart distribution

The Wishart distribution  $\text{Wis}(\mathbf{S}, k)$  is a multivariate generalisation of the gamma distribution  $\text{Gam}(\alpha, \theta)$  (Section 4.4) from one to  $d$  dimensions.

### Standard parametrisation

If the symmetric random matrix  $\mathbf{X}$  of dimension  $d \times d$  is Wishart distributed we write

$$\mathbf{X} \sim \text{Wis}(\mathbf{S}, k)$$

where  $\mathbf{S} = (s_{ij})$  is the scale parameter (a symmetric  $d \times d$  positive definite matrix with elements  $s_{ij}$ ). The dimension  $d$  is implicit in the scale parameter  $\mathbf{S}$ .

The shape parameter  $k$  takes on real values in the range  $k > d - 1$  and integer values in the range  $k \in 1, \dots, d - 1$  for  $d > 1$ . For  $k > d - 1$  the matrix  $\mathbf{X}$  is positive definite and invertible (see also Section 5.5), otherwise  $\mathbf{X}$  is singular and positive semi-definite.

The distribution has mean

$$\mathbb{E}(\mathbf{X}) = k\mathbf{S}$$

and variances of the elements of  $\mathbf{X}$  are

$$\text{Var}(x_{ij}) = k \left( s_{ij}^2 + s_{ii}s_{jj} \right)$$

The pdf is (for  $k > d - 1$ )

$$p(\mathbf{X}|\mathbf{S}, k) = \frac{1}{\Gamma_d(k/2) \det(2\mathbf{S})^{k/2}} \det(\mathbf{X})^{(k-d-1)/2} \exp(-\text{Tr}(\mathbf{S}^{-1}\mathbf{X})/2)$$

with the multivariate gamma function defined as

$$\Gamma_d(k/2) = \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma((k-j+1)/2)$$

Note that this pdf is a joint pdf over the  $d$  diagonal elements  $x_{ii}$  and the  $d(d-1)/2$  off-diagonal elements  $x_{ij}$  of the symmetric random matrix  $\mathbf{X}$ .

If  $\mathbf{S}$  is a scalar rather than a matrix (and hence  $d = 1$ ) then the multivariate Wishart distribution reduces to the univariate Wishart aka gamma distribution (Section 4.4).

The Wishart distribution is closely related to the multivariate normal distribution with mean zero. Specifically, if  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{S})$  then  $\mathbf{z}\mathbf{z}^T \sim \text{Wis}(\mathbf{S}, 1)$ .

#### R code

The `mnw` package implements the Wishart distribution. The pdf of the Wishart distribution is given by `mnw::dwish()`. The corresponding random number generator is `mnw::rwish()`.

### Mean parametrisation

It is useful to employ the Wishart distribution in **mean parametrisation**

$$\text{Wis}\left(\mathbf{S} = \frac{\mathbf{M}}{k}, k\right)$$

with parameters  $M = kS$  and  $k$ . In this parametrisation the mean is

$$E(X) = M = (\mu_{ij})$$

and variances of the elements of  $X$  are

$$\text{Var}(x_{ij}) = \frac{\mu_{ij}^2 + \mu_{ii}\mu_{jj}}{k}$$

### Special case: standard Wishart distribution

For  $S = I$  the Wishart distribution reduces to the **standard Wishart distribution**

$$X \sim \text{Wis}(I, k)$$

with a single shape parameter  $k$ . The mean is

$$E(X) = kI$$

and variances of the elements of  $X$  are

$$\text{Var}(x_{ij}) = \begin{cases} 2k & \text{if } i = j \\ k & \text{if } i \neq j \end{cases}$$

The pdf is (for  $k > d - 1$ )

$$p(X|k) = \frac{1}{\Gamma_d(k/2)2^{dk/2}} \det(X)^{(k-d-1)/2} \exp(-\text{Tr}(X)/2)$$

The standard Wishart distribution is closely related to the standard multivariate normal distribution with mean zero. Specifically, if  $z \sim N(0, I)$  then  $zz^T \sim \text{Wis}(I, 1)$ .

The **Bartlett decomposition** of the standard multivariate Wishart  $\text{Wis}(I, k)$  distribution for any real  $k > d - 1$  is obtained by Cholesky factorisation of the random matrix  $X = ZZ^T$ . By construction  $Z$  is a lower-triangular matrix with positive diagonal elements  $z_{ii}$  and lower off-diagonal elements  $z_{ij}$  with  $i > j$  and  $i, j \in \{1, \dots, d\}$ . The corresponding upper off-diagonal elements are set to zero ( $z_{ji} = 0$ ).

The  $d(d+1)/2$  elements of  $Z$  are independent and allow to generate a standard Wishart variate as follows:

- 1) the *squared* diagonal elements follow a univariate standard Wishart distribution  $z_{ii}^2 \sim \text{Wis}(1, k - i + 1)$  and
- 2) the off-diagonal elements follow the univariate standard normal distribution  $z_{ij} \sim N(0, 1)$ .
- 3) Then  $X = ZZ^T \sim \text{Wis}(I, k)$ .

### Scale transformation

If  $X \sim \text{Wis}(S, k)$  then the scaled symmetric random matrix  $AXA^T$  is also Wishart distributed with  $AXA^T \sim \text{Wis}(ASA^T, k)$  where the matrix  $A$  must be full rank and  $ASA^T$  remains positive definite. The matrix  $A$  may be rectangular, hence the size of  $AXA^T$  and  $ASA^T$  may be smaller compared to  $X$  and  $S$ .

The transformations between the Wishart distribution and the standard Wishart distribution are two important special cases:

- 1) With  $W^T W = S^{-1}$  and  $X \sim \text{Wis}(S, k)$  then  $Y = W X W^T \sim \text{Wis}(I, k)$  as  $W S W^T = I$ . This transformation reduces the Wishart distribution to the standard Wishart distribution.
- 2) Conversely, with  $LL^T = S$  and  $Y \sim \text{Wis}(I, k)$  then  $X = L Y L^T \sim \text{Wis}(S, k)$  as  $L I L^T = S$ . This transformation generates the Wishart distribution from the standard Wishart distribution.

### Convolution property

The convolution of  $n$  Wishart distributions with the same scale parameter  $S$  but possible different shape parameters  $k_i$  yields another Wishart distribution:

$$\sum_{i=1}^n \text{Wis}(S, k_i) \sim \text{Wis}\left(S, \sum_{i=1}^n k_i\right)$$

Note that the shape parameter  $k$  is restricted to be an integer in the range  $1, \dots, d - 1$  for  $d > 1$  but is a real number in the range  $k > d - 1$ . Thus, if the  $k_i$  are all valid shape parameters (for dimension  $d$ ) then  $\sum_{i=1}^n k_i$  is also a valid shape parameter.

Due the partial restriction of the shape parameter  $k$  to integer values the multivariate Wishart distribution is **not infinitely divisible** for  $d > 1$ .

The above includes the following construction of the multivariate Wishart distribution  $\text{Wis}(S, k)$  for integer-valued  $k$ . The sum of  $k$  independent Wishart random variables  $\text{Wis}(S, 1)$  with one degree of freedom and identical scale parameter yields a Wishart random variable  $\text{Wis}(S, k)$  with degree of freedom  $k$  and the same scale parameter. Thus, if  $z_1, z_2, \dots, z_k \sim N(0, S)$  are  $k$  independent samples from  $N(0, S)$  then  $\sum_{i=1}^k z_i z_i^T \sim \text{Wis}(S, k)$ .

## 5.5 Inverse Wishart distribution

The **inverse Wishart distribution**  $IW(\Psi, k)$  is a multivariate generalisation of the inverse gamma distribution  $IG(\alpha, \beta)$  (Section 4.5) from one to  $d$  dimensions. It is linked to the Wishart distribution  $Wis(S, k)$  (Section 5.4).

### Standard parametrisation

A symmetric positive definite random matrix  $X$  of dimension  $d \times d$  following an inverse Wishart distribution is denoted by

$$X \sim IW(\Psi, k)$$

where  $\Psi = (\psi_{ij})$  is the scale parameter (a  $d \times d$  positive definite symmetric matrix) and  $k > d - 1$  is the shape parameter. The dimension  $d$  is implicit in the scale parameter  $\Psi$ .

The mean is (for  $k > d + 1$ )

$$E(X) = \frac{\Psi}{k - d - 1}$$

and the variances of elements of  $X$  are (for  $k > d + 3$ )

$$\text{Var}(x_{ij}) = \frac{(k - d - 1)\psi_{ii}\psi_{jj} + (k - d + 1)\psi_{ij}^2}{(k - d)(k - d - 3)(k - d - 1)^2}$$

The inverse Wishart distribution  $IW(\Psi, k)$  has pdf

$$p(X|\Psi, k) = \frac{\det(\Psi/2)^{k/2}}{\Gamma_d(k/2)} \det(X)^{-(k+d+1)/2} \exp(-\text{Tr}(\Psi X^{-1})/2)$$

As with the Wishart distribution his pdf is a joint pdf over the  $d$  diagonal elements  $x_{ii}$  and the  $d(d-1)/2$  off-diagonal elements  $x_{ij}$  of the symmetric random matrix  $X$ .

The inverse Wishart and the Wishart distributions are linked. If  $X \sim IW(\Psi, k)$  then the inverse of  $X$  is Wishart distributed with inverted scale parameter:

$$X^{-1} \sim Wis(S = \Psi^{-1}, k)$$

where  $k$  is the shape parameter and  $S$  the scale parameter of the Wishart distribution.

If  $\Psi$  is a scalar  $\psi$  (and  $d = 1$ ) then the multivariate inverse Wishart distribution reduces to the univariate inverse Wishart distribution (Section 4.5).

#### R code

The `mniw` package implements the Wishart distribution. The pdf of the Wishart distribution is given by `mniw::diwish()`. The corresponding random number generator is `mniw::riwish()`.

### Mean parametrisation

Instead of  $\Psi$  and  $k$  we may also equivalently use  $M = \Psi/(k - d - 1)$  and  $\kappa = k - d - 1$  as parameters for the inverse Wishart distribution, so that

$$X \sim IW(\Psi = \kappa M, k = \kappa + d + 1)$$

with mean (for  $\kappa > 0$ )

$$E(X) = M$$

and variances (for  $\kappa > 2$ )

$$\text{Var}(x_{ij}) = \frac{\kappa \mu_{ii}\mu_{jj} + (\kappa + 2)\mu_{ij}^2}{(\kappa + 1)(\kappa - 2)}$$

For  $M$  equal to scalar  $\mu$  with  $d = 1$  the above reduces to the univariate inverse Wishart distribution in mean parametrisation.

### Biased mean parametrisation

Using  $T = (t_{ij}) = \Psi/(k - d + 1) = \Psi/\nu$  as biased mean parameter together with  $\nu = k - d + 1$  we arrive at the **biased mean parametrisation**

$$X \sim IW(\Psi = \nu T, k = \nu + d - 1)$$

The corresponding mean is (for  $\nu > 2$ )

$$E(X) = \frac{\nu}{\nu - 2} T = M$$

and the variances of elements of  $\mathbf{X}$  are (for  $\nu > 4$ )

$$\text{Var}(x_{ij}) = \left(\frac{\nu}{\nu-2}\right)^2 \frac{(\nu-2)t_{ii}t_{jj} + \nu t_{ij}^2}{(\nu-1)(\nu-4)}$$

As  $\mathbf{T} = \mathbf{M}(\nu-2)/\nu$  for large  $\nu$  the parameter  $\mathbf{T}$  will become identical to the true mean  $\mathbf{M}$ .

For  $\mathbf{T}$  equal to scalar  $\tau^2$  with  $d = 1$  the above reduces to the univariate inverse Wishart distribution in biased mean parametrisation.

### Scale transformation

If  $\mathbf{X} \sim \text{IW}(\boldsymbol{\Psi}, k)$  then the scaled symmetric random matrix  $\mathbf{A}\mathbf{X}\mathbf{A}^T$  is also inverse Wishart distributed with  $\mathbf{A}\mathbf{X}\mathbf{A}^T \sim \text{IW}(\mathbf{A}\boldsymbol{\Psi}\mathbf{A}^T, k)$  where the matrix  $\mathbf{A}$  has full rank and both  $\mathbf{A}\mathbf{X}\mathbf{A}^T$  and  $\mathbf{A}\boldsymbol{\Psi}\mathbf{A}^T$  remain positive definite. The matrix  $\mathbf{A}$  may be rectangular, hence the size of  $\mathbf{A}\mathbf{X}\mathbf{A}^T$  and  $\mathbf{A}\boldsymbol{\Psi}\mathbf{A}^T$  may be smaller compared to  $\mathbf{X}$  and  $\boldsymbol{\Psi}$ .

## 5.6 Multivariate $t$ -distribution

The **multivariate  $t$ -distribution**  $t_\nu(\boldsymbol{\mu}, \mathbf{T})$  is a multivariate generalisation of the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  (Section 4.6) from one to  $d$  dimensions. It is a generalisation of the multivariate normal distribution  $N(\boldsymbol{\mu}, \mathbf{T})$  (Section 5.3) with an additional parameter  $\nu > 0$  (degrees of freedom) controlling the probability mass in the tails.

Special cases include the **multivariate standard  $t$ -distribution**  $t_\nu(\mathbf{0}, \mathbf{I})$ , the **multivariate normal distribution**  $N(\boldsymbol{\mu}, \mathbf{T})$  and the **multivariate Cauchy** distribution  $\text{Cau}(\boldsymbol{\mu}, \mathbf{T})$ .

### Standard parametrisation

If  $\mathbf{x} \in \mathbb{R}^d$  is a multivariate  $t$ -distributed random variable we write

$$\mathbf{x} \sim t_\nu(\boldsymbol{\mu}, \mathbf{T})$$

where the vector  $\boldsymbol{\mu}$  is the location parameter (a  $d$  dimensional vector) and the dispersion parameter  $\mathbf{T}$  is a symmetric positive definite matrix of dimension  $d \times d$ . The dimension  $d$  is implicit in both parameters. The parameter  $\nu > 0$  prescribes the degrees of freedom. For small values of

$\nu$  the distribution is heavy-tailed and as a result only moments of order smaller than  $\nu$  are finite and defined.

The mean is (for  $\nu > 1$ )

$$\mathbb{E}(\mathbf{x}) = \boldsymbol{\mu}$$

and the variance (for  $\nu > 2$ )

$$\text{Var}(\mathbf{x}) = \frac{\nu}{\nu-2} \mathbf{T}$$

The pdf of  $t_\nu(\boldsymbol{\mu}, \mathbf{T})$  is

$$p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{T}, \nu) = \det(\mathbf{T})^{-1/2} \frac{\Gamma(\frac{\nu+d}{2})}{(\pi\nu)^{d/2} \Gamma(\frac{\nu}{2})} \left(1 + \frac{\Delta^2}{\nu}\right)^{-(\nu+d)/2}$$

with  $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{T}^{-1} (\mathbf{x} - \boldsymbol{\mu})$  the squared Mahalanobis distance between  $\mathbf{x}$  and  $\boldsymbol{\mu}$ . Note that this pdf is a joint pdf over the  $d$  elements  $x_1, \dots, x_d$  of the random vector  $\mathbf{x}$ .

For  $d = 1$  the random vector  $\mathbf{x} = x$  is a scalar,  $\boldsymbol{\mu} = \mu$ ,  $\mathbf{T} = \tau^2$  and thus the multivariate  $t$ -distribution reduces to the location-scale  $t$ -distribution (Section 4.6).

#### R code

The `mnormt` package implements the multivariate  $t$ -distribution. The function `mnormt::dmt()` provides the pdf and `mnormt::pmt()` returns the distribution function. The function `mnormt::rmt()` is the corresponding random number generator.

### Scale parametrisation

The multivariate  $t$ -distribution, like the multivariate distribution, can also be represented with a matrix scale parameter  $\mathbf{L}$  in place of a matrix dispersion parameter  $\mathbf{T}$ .

Let  $\mathbf{L}$  be a matrix scale parameter such that  $\mathbf{L}\mathbf{L}^T = \mathbf{T}$  and  $\mathbf{W} = \mathbf{L}^{-1}$  be the corresponding inverse matrix scale parameter with  $\mathbf{W}^T\mathbf{W} = \mathbf{T}^{-1}$ . By construction  $|\det(\mathbf{W})| = \det(\mathbf{T})^{-1/2}$  and  $|\det(\mathbf{L})| = \det(\mathbf{T})^{1/2}$ .

Note that  $\mathbf{T}$  alone does not fully determine  $\mathbf{L}$  and  $\mathbf{W}$  due to rotational freedom, see the discussion in Section 5.3 for details.

### Special case: multivariate standard $t$ -distribution

With  $\mu = \mathbf{0}$  and  $T = I$  the multivariate  $t$ -distribution reduces to the **multivariate standard  $t$ -distribution**  $t_\nu(\mathbf{0}, I)$ . It is a generalisation of the multivariate standard normal distribution  $N(\mathbf{0}, I)$  to allow for heavy tails.

The distribution has mean  $E(x) = \mathbf{0}$  (for  $\nu > 1$ ) and variance  $\text{Var}(x) = \frac{\nu}{\nu-2}I$  (for  $\nu > 2$ ).

The pdf of  $t_\nu(\mathbf{0}, I)$  is

$$p(x|\nu) = \frac{\Gamma(\frac{\nu+d}{2})}{(\pi\nu)^{d/2} \Gamma(\frac{\nu}{2})} \left(1 + \frac{x^T x}{\nu}\right)^{-(\nu+d)/2}$$

with the squared Mahalanobis distance reducing to  $\Delta^2 = x^T x$ .

For scalar  $x$  (and hence  $d = 1$ ) the multivariate standard  $t$ -distribution reduces to the Student's  $t$ -distribution  $t_\nu = t_\nu(0, 1)$ .

Unlike the multivariate standard normal distribution, the density of the multivariate standard  $t$ -distribution cannot be written as product of corresponding univariate standard densities.

### Special case: multivariate normal distribution

For  $\nu \rightarrow \infty$  the multivariate  $t$ -distribution  $t_\nu(\mu, T)$  reduces to the **multivariate normal distribution**  $N(\mu, T)$  (Section 5.3). Correspondingly, for  $\nu \rightarrow \infty$  the multivariate standard  $t$ -distribution  $t_\nu(\mathbf{0}, I)$  becomes equal to the **multivariate standard normal distribution**  $N(\mathbf{0}, I)$ .

This can be seen from the corresponding limits of the two factors in the pdf of the multivariate  $t$ -distribution that depend on  $\nu$ :

- 1) Following Sterling's approximation for large  $x$  we can approximate  $\log \Gamma(x) \approx (x-1)\log(x-1)$ . For large  $\nu$  this implies that

$$\frac{\Gamma((\nu+d)/2)}{(\pi\nu)^{d/2} \Gamma(\nu/2)} \rightarrow (2\pi)^{-d/2}$$

- 2) For small  $x$  we can approximate  $\log(1+x) \approx x$ . Thus for large  $\nu \gg d$  (and hence small  $\Delta^2/\nu$ ) this yields  $(\nu+d)\log(1+\Delta^2/\nu) \rightarrow \Delta^2$  and hence  $(1+\Delta^2/\nu)^{-(\nu+d)/2} \rightarrow e^{-\Delta^2/2}$ .

Hence, the pdf of  $t_\infty(\mu, T)$  is the multivariate normal pdf

$$p(x|\mu, T, \nu = \infty) = \det(T)^{-1/2} (2\pi)^{-d/2} e^{-\Delta^2/2}$$

### Special case: multivariate Cauchy distribution

For  $\nu = 1$  the multivariate  $t$ -distribution becomes the **multivariate Cauchy distribution**  $\text{Cau}(\mu, T) = t_1(\mu, T)$ .

Its mean, variance and other higher moments are all undefined.

It has pdf

$$p(x|\mu, T) = \det(T)^{-1/2} \Gamma\left(\frac{d+1}{2}\right) (\pi(1+\Delta^2))^{-(d+1)/2}$$

For scalar  $x$  (and hence  $d = 1$ ) the multivariate Cauchy distribution  $\text{Cau}(\mu, T)$  reduces to the univariate Cauchy distribution  $\text{Cau}(\mu, \tau^2)$ .

### Special case: multivariate standard Cauchy distribution

The **multivariate standard Cauchy distribution**  $\text{Cau}(\mathbf{0}, I) = t_1(\mathbf{0}, I)$  is obtained by setting  $\mu = \mathbf{0}$  and  $T = I$  in the multivariate Cauchy distribution or, equivalently, by setting  $\nu = 1$  in the multivariate standard  $t$ -distribution.

It has pdf

$$p(x) = \Gamma\left(\frac{d+1}{2}\right) (\pi(1+x^T x))^{-(d+1)/2}$$

For scalar  $x$  (and hence  $d = 1$ ) the multivariate standard Cauchy distribution  $\text{Cau}(\mathbf{0}, I)$  reduces to the standard univariate Cauchy distribution  $\text{Cau}(0, 1)$ .

### Location-scale transformation

Let  $L$  be a scale matrix for  $T$  and  $W$  the corresponding inverse scale matrix.

If  $x \sim t_\nu(\mu, T)$  then  $y = W(x - \mu) \sim t_\nu(\mathbf{0}, I)$ . This location-scale transformation reduces a multivariate  $t$ -distributed random variable to a standard multivariate  $t$ -distributed random variable.

Conversely, if  $y \sim t_\nu(\mathbf{0}, I)$  then  $x = \mu + Ly \sim t_\nu(\mu, T)$ . This location-scale transformation generates the multivariate  $t$ -distribution from the multivariate standard  $t$ -distribution.

Note that for  $\nu > 2$  under the location-scale transformation  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{L}\mathbf{y}$  with  $\text{Var}(\mathbf{y}) = \nu/(\nu - 2)\mathbf{I}$  we get  $\text{Cov}(\mathbf{x}, \mathbf{y}) = \nu/(\nu - 2)\mathbf{L}$ . This provides a means to choose between different factorisations of  $\mathbf{T}$  and  $\mathbf{T}^{-1}$ . For example, if positive correlation between corresponding elements in  $\mathbf{x}$  and  $\mathbf{y}$  is desired then the diagonal elements in  $\mathbf{L}$  must be positive.

For the special case of the multivariate Cauchy distribution (corresponding to  $\nu = 1$ ) similar relations hold between it and the multivariate standard Cauchy distribution. If  $\mathbf{x} \sim \text{Cau}(\boldsymbol{\mu}, \mathbf{T})$  then  $\mathbf{y} = \mathbf{W}(\mathbf{x} - \boldsymbol{\mu}) \sim \text{Cau}(\mathbf{0}, \mathbf{I})$ . Conversely, if  $\mathbf{y} \sim \text{Cau}(\mathbf{0}, \mathbf{I})$  then  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{L}\mathbf{y} \sim \text{Cau}(\boldsymbol{\mu}, \mathbf{T})$ .

### Convolution property

The multivariate  $t$ -distribution is not generally closed under convolution, with the exception of two special cases, the multivariate normal distribution ( $\nu = \infty$ ), see Section 5.3, and the multivariate Cauchy distribution ( $\nu = 1$ ) with the additional restriction that the dispersion parameters are proportional.

For the Cauchy distribution with  $\mathbf{T}_i = a_i^2 \mathbf{T}$ , where  $a_i > 0$  are positive scalars,

$$\sum_{i=1}^n \text{Cau}(\boldsymbol{\mu}_i, a_i^2 \mathbf{T}) \sim \text{Cau}\left(\sum_{i=1}^n \boldsymbol{\mu}_i, \left(\sum_{i=1}^n a_i\right)^2 \mathbf{T}\right)$$

### Multivariate $t$ -distribution as compound distribution

The multivariate  $t$ -distribution can be obtained as mixture of multivariate normal distributions with identical mean and varying covariance matrix. Specifically, let  $z$  be a univariate inverse Wishart random variable

$$z \sim \text{IW}(\psi = \nu, k = \nu) = \text{IG}\left(\alpha = \frac{\nu}{2}, \beta = \frac{\nu}{2}\right)$$

and let  $\mathbf{x}|z$  be multivariate normal

$$\mathbf{x}|z \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma} = z\mathbf{T})$$

The resulting marginal (scale mixture) distribution for  $\mathbf{x}$  is the multivariate  $t$ -distribution

$$\mathbf{x} \sim t_\nu(\boldsymbol{\mu}, \mathbf{T})$$

An alternative way to arrive at  $t_\nu(\boldsymbol{\mu}, \mathbf{T})$  is to include  $\mathbf{T}$  as parameter in the inverse Wishart distribution

$$\mathbf{Z} \sim \text{IW}(\boldsymbol{\Psi} = \nu\mathbf{T}, k = \nu + d - 1)$$

and let

$$\mathbf{x}|\mathbf{Z} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma} = \mathbf{Z})$$

Note that  $\mathbf{T}$  is now the biased mean parameter of the multivariate inverse Wishart distribution. This characterisation is useful in Bayesian analysis.



## Bibliography

- Mardia, K. V., J. T. Kent, and J. M. Bibby. 1979. *Multivariate Analysis*. Academic Press.
- Whittle, P. 2000. *Probability via Expectation*. 3rd ed. Springer. <https://doi.org/10.1007/978-1-4612-0509-8>.