

# **Probability and Distribution Refresher**

Korbinian Strimmer

28 July 2024

# 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 . . . . .	3
1.3 De Moivre-Sterling approximation of the factorial . . . . .	4
1.4 Multinomial and binomial coefficient . . . . .	4
<b>2 Probability</b>	<b>6</b>
2.1 Random variables . . . . .	6
2.2 Probability mass and density function . . . . .	7
2.3 Distribution function and quantile function . . . . .	8
2.4 Families of distributions . . . . .	9
2.5 Expectation of a random variable . . . . .	9
2.6 Jensen's inequality for the expectation . . . . .	10
2.7 Probability as expectation . . . . .	10
2.8 Moments and variance of a random variable . . . . .	10
2.9 Random vectors and their mean and variance . . . . .	11
2.10 Correlation matrix . . . . .	12
<b>3 Transformations</b>	<b>13</b>
3.1 Affine or location-scale transformation of random variables	13
3.2 General invertible transformation of random variables . .	15
3.3 Exponential tilting and exponential families . . . . .	16
3.4 Sums of random variables and convolution . . . . .	17
<b>4 Univariate distributions</b>	<b>20</b>
4.1 Binomial distribution . . . . .	20
4.2 Beta distribution . . . . .	23
4.3 Normal distribution . . . . .	26

4.4	Gamma distribution . . . . .	28
4.5	Inverse gamma distribution . . . . .	34
4.6	Location-scale $t$ -distribution . . . . .	39
<b>5</b>	<b>Multivariate distributions</b>	<b>44</b>
5.1	Multinomial distribution . . . . .	44
5.2	Dirichlet distribution . . . . .	47
5.3	Multivariate normal distribution . . . . .	50
5.4	Wishart distribution . . . . .	54
5.5	Inverse Wishart distribution . . . . .	58
5.6	Multivariate $t$ -distribution . . . . .	60
	<b>Bibliography</b>	<b>66</b>



# Welcome

The Probability and Distribution Refresher notes were written by [Korbinian Strimmer](#) from 2018–2024. This version is from 28 July 2024.

If you have any questions, comments, or corrections please get in touch!

<sup>1</sup>

## Updates

The notes will be updated from time to time. To view the current version visit the

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

You may also wish to 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.strimmer@manchester.ac.uk](mailto:korbinian.strimmer@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

Summation:

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

Multiplication:

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

Indicator function (in Iverson bracket notation):

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

Scalar: plain type, typically lower case ( $x$ ,  $\theta$ ), sometimes upper case ( $K$ ).

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

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

## 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**

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

## 1 Combinatorics

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

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



#### 1.4 Multinomial and binomial coefficient

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 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** or **state space** of the random variable and is denoted by  $\Omega = \{\omega_1, \omega_2, \dots\}$ . The outcomes  $\omega_i$  are the **elementary events**. The sample space  $\Omega$  can be finite or infinite. Depending on type of outcomes the random variable is **discrete** or **continuous**.

An event  $A \subseteq \Omega$  is a subset of  $\Omega$  and thus itself a set composed of elementary events:  $A = \{a_1, a_2, \dots\}$ . This includes as special cases the full set  $A = \Omega$ , the empty set  $A = \emptyset$ , and the elementary events  $A = \omega_i$ . The complementary event  $A^C$  is the complement of the set  $A$  in the set  $\Omega$  so that  $A^C = \Omega \setminus A = \{\omega_i \in \Omega : \omega_i \notin A\}$ .

The probability of an event  $A$  is denoted by  $\Pr(A)$ . Essentially, to obtain this probability we need to count the elementary elements corresponding to  $A$ . To do this we assume as **axioms of probability** that

- $\Pr(A) \geq 0$ , probabilities are positive,
- $\Pr(\Omega) = 1$ , the certain event has probability 1, and
- $\Pr(A) = \sum_{a_i \in A} \Pr(a_i)$ , the probability of an event equals the sum of its constituting elementary events  $a_i$ . This sum is taken over a finite or countable infinite number of elements.

This implies

- $\Pr(A) \leq 1$ , i.e. probabilities all lie in the interval  $[0, 1]$
- $\Pr(A^C) = 1 - \Pr(A)$ , and
- $\Pr(\emptyset) = 0$

Assume now that we have two events  $A$  and  $B$ . The probability of the event “ $A$  and  $B$ ” is then given by the probability of the set intersection  $\Pr(A \cap B)$ . Likewise the probability of the event “ $A$  or  $B$ ” is given by the probability of the set union  $\Pr(A \cup B)$ .

From the above it is clear that the definition and theory of probability is closely linked to set theory, and in particular to measure theory. Indeed, viewing probability as a special type of measure allows for an elegant treatment of both discrete and continuous random variables.

## 2.2 Probability mass and density function

To describe a random variable  $x$  with state space  $\Omega$  we need a way to effectively store the probabilities of the corresponding elementary outcomes  $x \in \Omega$ .

**For simplicity of notation we use the same symbol to denote the random variable and its elementary outcomes.**<sup>1</sup> This convention greatly facilitates working with random vectors and matrices and follows, e.g., the classic multivariate statistics textbook by Mardia, Kent, and Bibby (1979). If a quantity is random we will always specify this explicitly in the context.

For a discrete random variable we define the event  $A = \{x : x = a\} = \{a\}$  and get the probability

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

directly from the **probability mass function** (pmf), here denoted by lower case  $f$  (but we frequently also use  $p$  or  $q$ ). The pmf has the property that  $\sum_{x \in \Omega} f(x) = 1$  and that  $f(x) \in [0, 1]$ .

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

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

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.

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

---

<sup>1</sup>For scalar random variables many texts use upper case to designate the random variable and lower case for its realisations. However, this convention quickly breaks down in multivariate statistics when dealing with random vectors and random matrices. Hence, we use upper case primarily to indicate a matrix quantity (in bold type). Upper case (in plain type) may denote sets and some scalar quantities traditionally written in upper case (e.g.  $R^2$ ,  $K$ ).

## 2 Probability

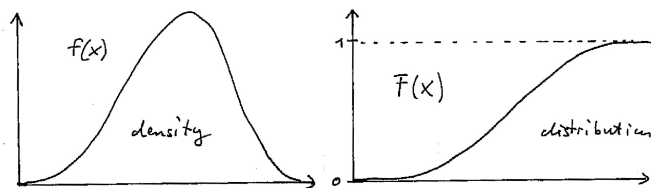


Figure 2.1: Density function and distribution function.

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

### 2.3 Distribution function and quantile function

As alternative to using the pdmf we may use a **distribution function** to describe the random variable. This assumes that an ordering exist among the elementary events so that we can define the event  $A = \{x : x \leq a\}$  and compute its probability as

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

This is also known **cumulative distribution function** (cdf) and is denoted by upper case  $F$  (or  $P$  and  $Q$ ). By construction the distribution function is monotonically non-decreasing and its value ranges from 0 to 1. With its help we can compute the probability of an interval set such as

$$\Pr(a < x \leq b) = F(b) - F(a).$$

The inverse of the distribution function  $y = F(x)$  is the **quantile function**  $x = F^{-1}(y)$ . The 50% quantile  $F^{-1}(\frac{1}{2})$  is called the **median**.

If the random variable  $x$  has distribution function  $F$  we write  $x \sim F$ .

Figure 2.1 illustrates a density function  $f(x)$  and the corresponding distribution function  $F(x)$ .

## 2.4 Families of distributions

A distribution  $F_\theta$  with a parameter  $\theta$  constitutes a **distribution family** collecting all the distributions corresponding to particular instances of the parameter. The parameter  $\theta$  therefore acts as an index of the distributions contained in the family.

The corresponding pmf is written either as  $f_\theta(x)$ ,  $f(x; \theta)$  or  $f(x|\theta)$ . The latter form is the most general as it suggests that the parameter  $\theta$  may potentially also have its own distribution, with a joint density formed by  $f(x, \theta) = f(x|\theta)f(\theta)$ .

Note that any parametrisation is generally not unique, as a one-to-one transformation of  $\theta$  will yield another equivalent index to the same distribution family. Typically, for most commonly used distribution families there are several standard parametrisations. Often we use those parametrisations where the parameters can be interpreted easily (e.g. in terms of moments).

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  $\theta$ .

## 2.5 Expectation of a random variable

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

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

Note the notation to emphasise that the expectation is taken with regard to the distribution  $F$ . The subscript  $F$  is usually left out if there are no ambiguities. Furthermore, because the sum or integral may diverge the expectation is not necessarily always defined (in contrast to quantiles).

The expected value of a function of a random variable  $h(x)$  is obtained similarly:

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

## 2 Probability

This is called the “[law of the unconscious statistician](#)”, or short LOTUS. Again, to highlight that the random variable  $x$  has distribution  $F$  we write  $E_F(h(x))$ .

### 2.6 Jensen’s inequality for the expectation

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

$$E(h(x)) \geq h(E(x))$$

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

### 2.7 Probability as expectation

Probability itself can also be understood as an expectation. For an event  $A$  we can define a corresponding indicator function  $[x \in A]$  for an elementary element  $x$  to be part of  $A$ . From the above it then follows

$$E([x \in A]) = \Pr(A)$$

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](#)).

### 2.8 Moments and variance of a random variable

The moments of a random variable are defined as follows:

- Zeroth moment:  $E(x^0) = 1$  by construction of a pmf,
- First moment:  $E(x^1) = E(x) = \mu$ , the mean,
- Second moment:  $E(x^2)$
- The variance is the second moment centred about the mean  $\mu$ :

$$\text{Var}(x) = E\left((x - \mu)^2\right) = \sigma^2$$

- The variance can also be computed by  $\text{Var}(x) = E(x^2) - E(x)^2$ . This provides an example of Jensen's inequality, with  $E(x^2) = E(x)^2 + \text{Var}(x) \geq E(x)^2$ .

A distribution does not necessarily need to have any finite first or higher moments. An example is the location-scale  $t$ -distribution (Section 4.6) that depending on the value of the parameter  $\nu$  may not have a mean or variance (or other higher moments).

## 2.9 Random vectors and their mean and variance

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

For a random vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \sim F$  the mean  $E(\mathbf{x}) = \boldsymbol{\mu}$  is given by the means of its elements, i.e.  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^T$  with  $\mu_i = E(x_i)$ . Thus, the mean of a random vector of dimension  $d$  is a vector of the same length.

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

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

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

<sup>2</sup>In our notational conventions, a **vector**  $\mathbf{x}$  is written in *lower case in bold type*, a **matrix**  $\mathbf{M}$  in *upper case in bold type*. Hence random vectors and matrices as well as their realisations are indicated in bold type, with vectors given in lower case and matrices in upper case. Hence, as for scalar variables, upper vs. lower case does not indicate randomness vs. realisation.

## 2 Probability

$\sigma_{ii}$  correspond to the variances  $\text{Var}(x_i) = \sigma_i^2$  of the elements of  $\mathbf{x}$ . The covariance matrix is by construction **positive semi-definite**, i.e. the eigenvalues of  $\mathbf{\Sigma}$  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.

### 2.10 Correlation matrix

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

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

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

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

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

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

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



# 3 Transformations

## 3.1 Affine or location-scale transformation of random variables

### Affine transformation

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

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

### 3 Transformations

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

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

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

$$E(\mathbf{y}) = \mathbf{a} + \mathbf{B} \boldsymbol{\mu}_x$$

and

$$\text{Var}(\mathbf{y}) = \mathbf{B} \boldsymbol{\Sigma}_x \mathbf{B}^T$$

where  $E(\mathbf{x}) = \boldsymbol{\mu}_x$  and  $\text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}_x$  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 of random variables

### General invertible transformation

As above we assume  $x \sim F_x$  is a scalar random variable and  $\mathbf{x} \sim F_{\mathbf{x}}$  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  $E(\mathbf{x}) = \boldsymbol{\mu}_x$  and  $\text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}_x$  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  $\boldsymbol{\mu}_x$  and the mean and variance resulting from the linear transformation is reported.

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

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

### 3 Transformations

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  $u(x)$  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  $u(x)$  is small.

The exponential tilt of a **base distribution**  $P_0$  with pdmf  $p_0(x)$  towards the linear combination  $\eta^T u(x)$  of the canonical statistics  $u(x)$  and the

**canonical parameters**  $\eta$  yields the distribution family  $P_\eta$  with pdmf

$$\begin{aligned} p(x|\eta) &= e^{\eta^T u(x)} b(x) / e^{\psi(\eta)} \\ &= \underbrace{e^{\eta^T u(x)}}_{\text{exponential tilt}} p_0(x) / e^{\psi(\eta) - \psi(0)} \end{aligned}$$

where  $b(x)$  is a positive base function. The normalising factor  $e^{\psi(\eta)}$  ensures that  $p(x|\eta)$  integrates to one. The pdmf of the base distribution is given by  $p_0(x) = b(x)/e^{\psi(0)}$ .

The distribution family  $P_\eta$  obtained by exponential tiling is called an **exponential family**. The corresponding log-pdmf is

$$\log p(x|\eta) = \eta^T u(x) + \log b(x) - \psi(\eta)$$

The **log-normaliser** or **log-partition function**  $\psi(\eta)$  is obtained by computing

$$\psi(\eta) = \log \int_x e^{\eta^T u(x)} b(x) dx$$

The set of values of  $\eta$  for which the integral is finite and hence for which  $\psi(\eta) < \infty$  defines the parameter space of the exponential family. Some choices of  $b(x)$  and  $u(x)$  will not allow for a finite normalising factor for any  $\eta$  and hence these cannot be used to form an exponential family.

Many commonly used distribution families are exponential families (most importantly the normal distribution). Exponential families are extremely important in probability and statistics. They provide highly effective models for statistical learning using entropy, likelihood and Bayesian approaches, allow for substantial data reduction via minimal sufficiency, and provide the basis of generalised linear models. Furthermore, exponential families often enable to generalise probabilistic results valid for the normal distribution to more general settings.

## 3.4 Sums of random variables and convolution

### Moments

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

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

### 3 Transformations

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 repeated application of **convolution** (symbolised by the  $*$  operator):

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

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 it can be applied in any order to compute the convolution of multiple functions. Furthermore, the convolution of pdfms yields another pdmf.

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: [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.

## 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

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

and the variance is

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



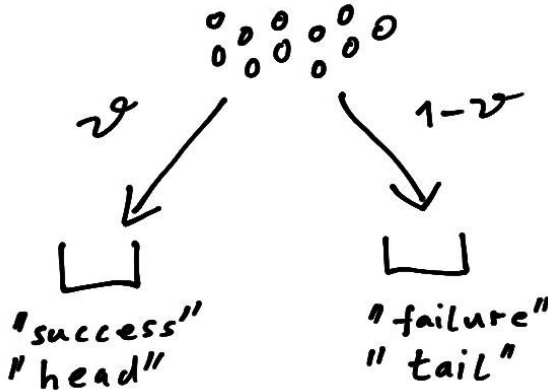


Figure 4.1: Binomial urn model.

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 family 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.

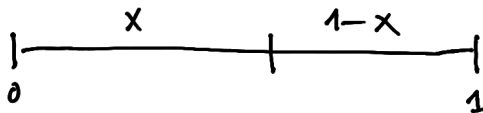


Figure 4.2: Stick breaking visualisation of a beta random variable.

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).

#### 4 Univariate distributions

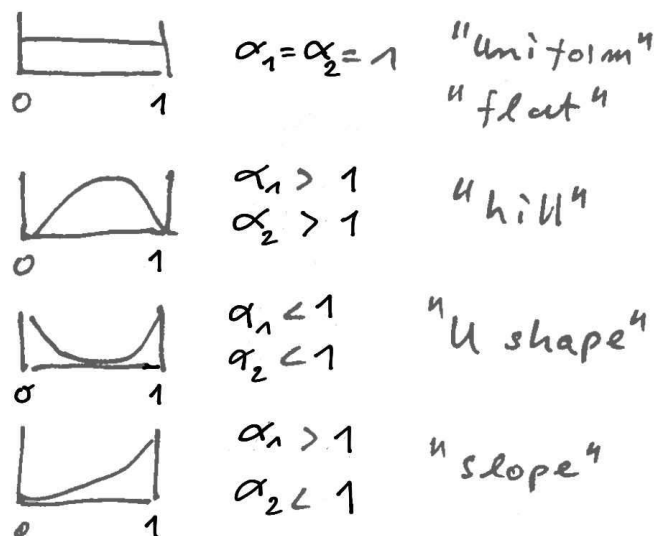


Figure 4.3: Shapes of the pdf of the beta distribution.

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

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

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 \(1777–1855\)](#).

A special case is the **standard normal distribution**  $N(0, 1)$ .

### 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()`.

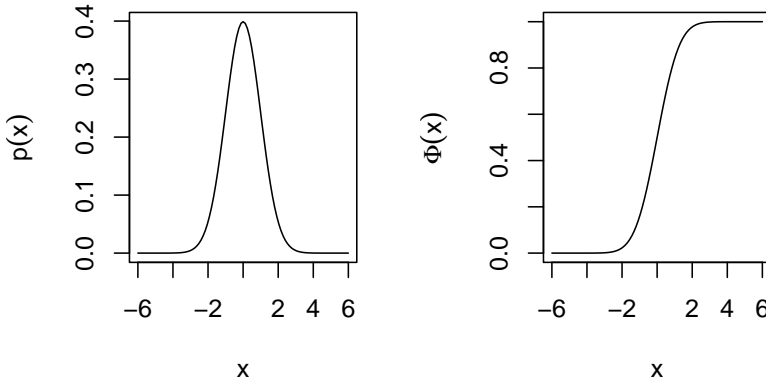


Figure 4.4: Probability density function (left) and cumulative density function (right) of the standard normal distribution.

### 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.

### 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**.

## 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



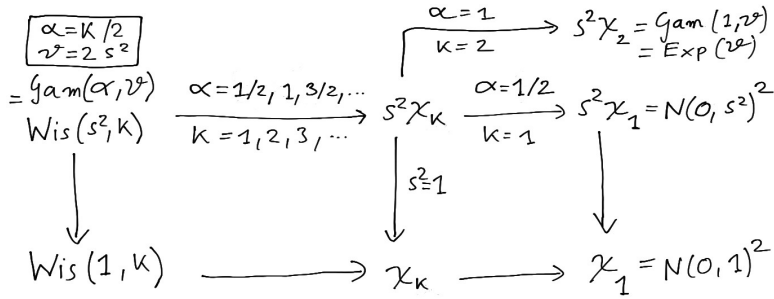


Figure 4.5: The gamma and the univariate Wishart distribution and its relatives.

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 quan-

tile 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  $\alpha \in \{1/2, 1, 3/2, 2, \dots\}$ .

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

$$E(x) = k$$

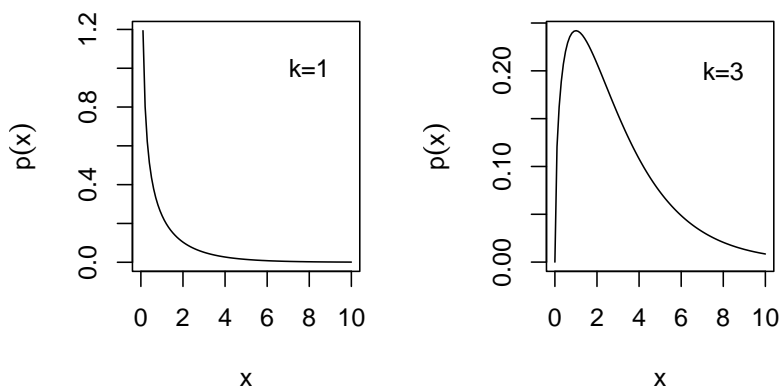


Figure 4.6: Probability density function of the chi-squared distribution.

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:

- a) 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.
- b) 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**

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

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  $k > 2$ )

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

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  $x \sim IW(\psi, k)$  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  $\nu > 2$ )

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

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  $IW(\psi, k)$  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**

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

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.

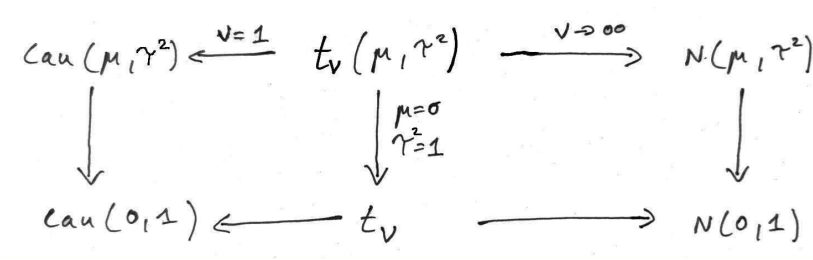
### Scale transformation

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  $1/x \sim \psi^{-1} \chi_k^2$ .

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

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  $\text{Inv-}\chi^2(\nu, \tau^2)$ . 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  $\nu > 1$ )

$$E(x) = \mu$$

## 4 Univariate distributions

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

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

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, `pcauchy()` is the distribution function and `qcauchy()` the quantile function. The corresponding random number generator is `rcauchy()`.

### 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  $\mathbf{x}$  describes the allocation of  $n$  items to  $K$  classes. We write

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

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  $\mathbf{x} = (x_1, \dots, x_K)^T$  is an integer  $x_i \in \{0, 1, \dots, n\}$  and  $\mathbf{x}$  satisfies the constraint  $\mathbf{x}^T \mathbf{1}_K = \sum_{i=1}^K x_i = n$ . Therefore the support of  $\mathbf{x}$  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

$$\mathbb{E}(\mathbf{x}) = n\theta$$

The covariance matrix is

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



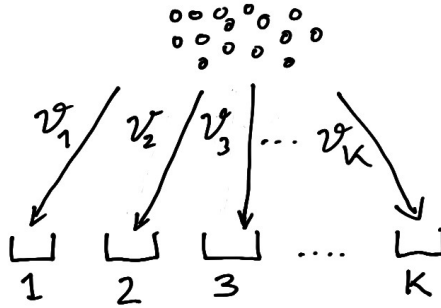


Figure 5.1: Multinomial urn model.

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

The corresponding pmf is

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

where  $\binom{n}{x_1, \dots, x_n}$  is the multinomial coefficient. While all  $K$  elements of  $x$  appear in the pmf recall that due 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

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

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

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  $\boldsymbol{\theta}$  but possibly different number of items  $n_i$ , yields another multinomial distribution with the same parameter  $\boldsymbol{\theta}$ :

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

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

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

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**  $\text{Dir}(\boldsymbol{\alpha})$  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.

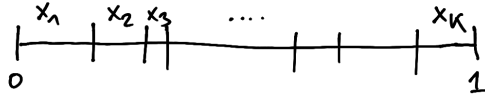


Figure 5.2: Stick breaking visualisation of a Dirichlet random variable.

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

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

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  $\alpha$  defined as

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

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

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

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

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

$$E(x) = \mu$$

and

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

The covariance matrix is singular by construction because of the dependencies among the elements of  $\mathbf{x}$ . 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  $p(\mathbf{x}) = 1/\Gamma(K)$ . 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.

A special case is the **multivariate standard normal distribution**  $N(\mathbf{0}, I)$ .

### 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  $\mathbf{x} \in \mathbb{R}^d$ .

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

$$\mathbf{x} \sim N(\mu, \Sigma)$$

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  $\boldsymbol{\mu} = \mathbf{0}$  and variance  $\boldsymbol{\Sigma} = \mathbf{I}$ . 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.

### Location-scale transformation

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

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

Note that under the location-scale transformation  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{L}\mathbf{y}$  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  $\boldsymbol{\Sigma}$  and  $\boldsymbol{\Sigma}^{-1}$ . 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:

$$\sum_{i=1}^n N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \sim N\left(\sum_{i=1}^n \boldsymbol{\mu}_i, \sum_{i=1}^n \boldsymbol{\Sigma}_i\right)$$

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  $\text{Wis}(S, k)$  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  $X$  of dimension  $d \times d$  is Wishart distributed we write

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

where  $S = (s_{ij})$  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  $S$ .

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  $X$  is positive definite and invertible (see also Section 5.5), otherwise  $X$  is singular and positive semi-definite.

The distribution has mean

$$E(X) = kS$$

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  $k > d - 1$ )

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

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 `mniw` package implements the Wishart distribution. The pdf of the Wishart distribution is given by `mniw::dwish()`. The corresponding random number generator is `mniw::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  $\mathbf{M} = k\mathbf{S}$  and  $k$ . In this parametrisation the mean is

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

and variances of the elements of  $\mathbf{X}$  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  $z_1, z_2, \dots, z_k \sim N(\mathbf{0}, S)$  are  $k$  independent samples from  $N(\mathbf{0}, S)$  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

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

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 `mnw` package implements the Wishart distribution. The pdf of the Wishart distribution is given by `mnw::diwish()`. The corresponding random number generator is `mnw::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  $X$  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  $T = M(\nu - 2)/\nu$  for large  $\nu$  the parameter  $T$  will become identical to the true mean  $M$ .

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

## Scale transformation

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

## 5.6 Multivariate $t$ -distribution

The **multivariate  $t$ -distribution**  $t_\nu(\mu, T)$  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  $N(\mu, T)$  (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**  $t_\nu(0, I)$ , the **multivariate normal distribution**  $N(\mu, T)$  and the **multivariate Cauchy** distribution  $\text{Cau}(\mu, T)$ .

### Standard parametrisation

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

$$x \sim t_\nu(\mu, T)$$

where the vector  $\mu$  is the location parameter (a  $d$  dimensional vector) and the dispersion parameter  $T$  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  $\nu > 1$ )

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

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

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

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

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

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

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

### 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

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

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

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

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  $t_\nu(\boldsymbol{\mu}, T)$  is to include  $T$  as parameter in the inverse Wishart distribution

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

and let

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

Note that  $T$  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>.