

# **Probability and Distribution Refresher**

Korbinian Strimmer

18 December 2025

# Table of contents

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

<b>4 Evaluation</b>	<b>24</b>
4.1 Loss functions . . . . .	24
4.2 Common loss functions . . . . .	25
4.3 Scoring rules . . . . .	26
4.4 Common scoring rules . . . . .	29
<b>5 Univariate distributions</b>	<b>34</b>
5.1 Binomial distribution . . . . .	34
5.2 Beta distribution . . . . .	37
5.3 Normal distribution . . . . .	40
5.4 Gamma distribution . . . . .	43
5.5 Inverse gamma distribution . . . . .	49
5.6 Location-scale $t$ -distribution . . . . .	54
<b>6 Multivariate distributions</b>	<b>59</b>
6.1 Multinomial distribution . . . . .	59
6.2 Dirichlet distribution . . . . .	62
6.3 Multivariate normal distribution . . . . .	65
6.4 Wishart distribution . . . . .	69
6.5 Inverse Wishart distribution . . . . .	73
6.6 Multivariate $t$ -distribution . . . . .	76
<b>7 Exponential families</b>	<b>81</b>
7.1 Definition of an exponential family . . . . .	81
7.2 Roles of the partition function . . . . .	83
7.3 Further properties . . . . .	85
7.4 Univariate exponential families . . . . .	87
7.5 Multivariate exponential families . . . . .	89
<b>Bibliography</b>	<b>91</b>

# Welcome

The Probability and Distribution Refresher notes were written by [Korbinian Strimmer](#) from 2018–2025. This version is from 18 December 2025.

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

## Updates

The lecture notes will be updated from time to time.

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

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

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

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

## License

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

---

<sup>1</sup>Email address: [korbinian.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 provide a quick refresher on some essential combinatorics and probability concepts, and present an overview of selected univariate and multivariate distributions, along with an introduction to exponential families.

The notes are supporting information for a number of lecture courses in statistics that I teach or have taught at the [Department of Mathematics of the University of Manchester](#).

This includes the currently offered modules:

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

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

- [MATH20802 Statistical Methods](#).

# 1 Combinatorics

## 1.1 Some basic mathematical notation

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

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

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

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

Summation:

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

Product:

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

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

Indicator function (in Iverson bracket notation):

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

## 1.2 Number of permutations

A **permutation** or **ordering** is a specific arrangement of items in a sequence, or equivalently, a specific assignment to labelled positions.

The **factorial**

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

is the number of permutations of  $n$  distinct items, where  $n$  is a positive integer.

For  $n = 0$  the factorial is defined as

$$0! = 1$$

Thus, the factorial  $n!$  equals the number of ways to place  $n$  distinct items into  $n$  labelled boxes so that each box contains exactly one item.

## 1.3 Multinomial and binomial coefficient

The **multinomial coefficient** for  $K$  groups

$$\begin{aligned} W_K &= \binom{n}{n_1, \dots, n_K} \\ &= \frac{n!}{n_1! n_2! \dots n_K!} \end{aligned}$$

is the number of permutations of  $n$  distinct items allocated to  $K$  groups, with  $n_k$  unordered items in group  $k$  and  $\sum_{k=1}^K n_k = n$ .

Thus, the multinomial coefficient  $W_K$  equals the number of ways to place  $n$  distinct items into  $K$  labelled boxes so that box  $k$  contains exactly  $n_k$  unordered items, with  $\sum_{k=1}^K n_k = n$ .

For  $n_k = 1$  (and thus  $K = n$ ) the multinomial coefficient reduces to the factorial.

For two groups ( $K = 2$ ) the multinomial coefficient becomes the **binomial coefficient**

$$\begin{aligned} W_2 &= \binom{n}{n_1, n_2} = \binom{n}{n_1, n - n_1} \\ &= \frac{n!}{n_1! (n - n_1)!} \\ &= \binom{n}{n_1} \end{aligned}$$

## 1.4 De Moivre-Sterling approximation

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

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

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

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

The de Moivre-Sterling approximation applied to the multinomial coefficient yields

$$\begin{aligned} \log W_K &\approx -n \sum_{k=1}^K \frac{n_k}{n} \log \left(\frac{n_k}{n}\right) \\ &= -n \sum_{k=1}^K q_k \log q_k = nH(\hat{Q}) \end{aligned}$$

Hence, for large  $n$  and large  $n_k$  the logarithm of the multinomial coefficient equals  $n$  times the information entropy  $H(\hat{Q})$  of the empirical categorical distribution  $\hat{Q}$  with class frequencies  $\hat{q}_k = n_k/n$ .

## 2 Probability

### 2.1 Random variables

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

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

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

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

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

This implies

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

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

### 2.2 Conditional probability

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

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

This identity follows from the axioms.

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

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

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

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

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

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

which in turn yields **Bayes' theorem**

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

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

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

## 2.3 Probability mass and density function

The **distribution** (or **law**) of a random variable  $x$  with sample space  $\Omega$  is the probability measure that assigns probabilities to values or ranges of  $x$ . This is done in practise by employing probability mass functions (pmf, for discrete random variables) or probability density functions (pdf, for continuous random variables).

The scalar random variable  $x$  is written in lowercase plain font. We use the same symbol  $x$  for both the random variable and its realisations.<sup>1</sup>

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

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

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

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

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

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

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

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

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

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

<sup>1</sup>This notation is common in statistical machine learning and multivariate statistics, see for example Mardia, Kent, and Bibby (1979). An alternative convention uses uppercase letters for random variables and lowercase for outcomes, but that convention is problematic for multivariate objects (random vectors and random matrices) and is also ill-suited in Bayesian statistics where parameters are modelled as random variables.

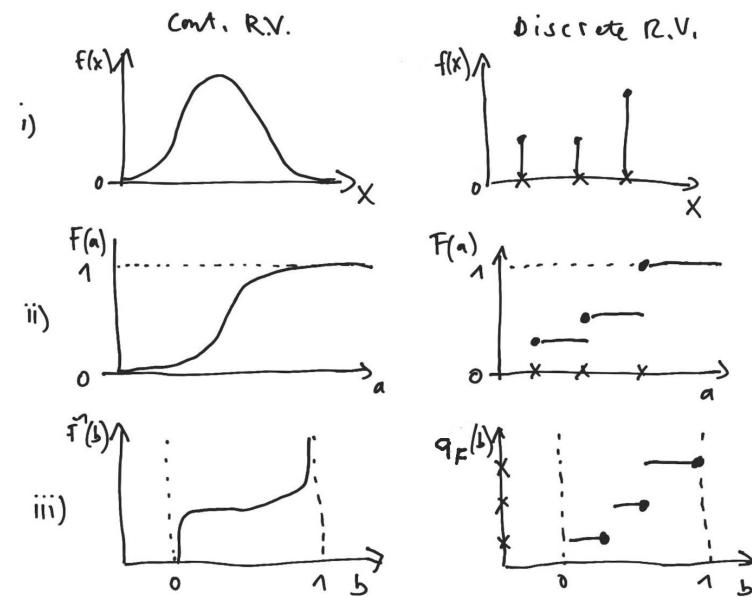


Figure 2.1: Illustration of i) pdmf, ii) distribution function and iii) quantile function for a continuous (first column) and a discrete random variable (second column).

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

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

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

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

## 2.4 Cumulative distribution function

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

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

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

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

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

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

This works both for discrete and continuous random variables.

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

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

## 2.5 Quantile function and quantiles

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

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

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

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

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

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

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

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

## 2.6 Expectation or mean

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

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

The subscript  $F$  in  $E_F(x)$  indicates that the expectation is taken with regard to the distribution  $F$ , but is usually left out if there are no ambiguities. The notation  $E(F)$  emphasises that the mean is a functional of the distribution  $F$ .

Because the sum or integral may diverge, not all distributions have finite means so the mean does not always exist (in contrast to the median, or quantiles in general). For example, the location-scale  $t$ -distribution  $t_v(\mu, \tau^2)$  does not have a mean for a degree of freedom in the range  $0 < v \leq 1$  (see Section 5.6).

Expectation is a **linear operator**, meaning that

$$E(a_1 x_1 + a_2 x_2) = a_1 E(x_1) + a_2 E(x_2)$$

for random variables  $x_1 \sim F_1$  and  $x_2 \sim F_2$  and constants  $a_1$  and  $a_2$ .

Consequently, expectation is **mixture preserving** so that

$$E(Q_\lambda) = (1 - \lambda) E(Q_0) + \lambda E(Q_1)$$

for the mixture  $Q_\lambda = (1 - \lambda)Q_0 + \lambda Q_1$  with  $0 < \lambda < 1$  and  $Q_0 \neq Q_1$ .

## 2.7 Variance

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

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

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

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

Like the mean, the variance may diverge and hence not necessarily exists for all distribution. For example, the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  does not have a variance for the degree of freedom in the range  $0 < \nu \leq 2$  (see Section 5.6).

## 2.8 Moments of a distribution

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

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

Special important cases are the

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

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

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

The first few central moments are the

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

The moments of a distribution are not necessarily all finite, i.e. some moments may not exist. For example, the location-scale  $t$ -distribution  $t_\nu(\mu, \tau^2)$  only has finite moments of degree smaller than the degree of freedom  $\nu$  (see Section 5.6).

## 2.9 Expectation of a transformed random variable

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

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

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

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

## 2.10 Probability as expectation

Probability itself can also be understood as an expectation.

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

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

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

## 2.11 Jensen's inequality for the expectation

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

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

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

An example of Jensen's inequality is  $\mathbb{E}(x^2) \geq \mathbb{E}(x)^2$ .

## 2.12 Random vectors and their mean and variance

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

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

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

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

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

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

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

The elements  $\text{Cov}(x_i, x_j) = \sigma_{ij}$  describe the covariance between the random variables  $x_i$  and  $x_j$ . The covariance matrix is symmetric, hence  $\sigma_{ij} = \sigma_{ji}$ . The diagonal elements  $\text{Cov}(x_i, x_i) = \sigma_{ii}$  correspond to the individual variances  $\text{Var}(x_i) = \sigma_i^2$ . By construction, the covariance matrix  $\Sigma$  is **positive semi-definite**, i.e. the eigenvalues of  $\Sigma$  are all positive or equal to zero.

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

## 2.13 Correlation matrix

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

Specifically, denote by  $V$  the diagonal matrix containing the variances

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

then the correlation matrix  $P$  is given by

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

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

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

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

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

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

## 2.14 Parameters and families of distributions

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

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

If a random variable  $x$  has distribution  $F(\theta)$  we write  $x \sim F(\theta)$ . In case of named distributions, we use the corresponding abbreviation such  $x \sim N(\mu, \sigma^2)$  (normal distribution) or  $x \sim \text{Beta}(\alpha_1, \alpha_2)$ .

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

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

If distinct parameters correspond to distinct distributions they are called **identifiable**. This is an important property as it allows parameters to be estimated from data. Specifically, if parameters are identifiable then  $P(\theta_1) = P(\theta_2)$  implies  $\theta_1 = \theta_2$ , and conversely if  $P(\theta_1) \neq P(\theta_2)$  then  $\theta_1 \neq \theta_2$ . If parameters and distributions are unique within a neighbourhood  $\theta_0 + \epsilon$  relative to a reference value  $\theta_0$ , rather than globally, they are **locally identifiable** at  $\theta_0$ .

An important class of distribution families are **exponential families** discussed in Chapter 7.

## 3 Transformations and convolution

### 3.1 Affine or location-scale transformation

#### Transformation rule

Suppose  $x$  is a scalar. The 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 vector  $x$  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  $a = 0$  this is a **linear transformation**.

For  $m = d$  (square  $B$ ) and  $\det(B) \neq 0$  the affine transformation is **invertible** with back-transformation

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

## Probability mass function

If  $x \sim F_x$  is a discrete scalar random variable with pmf  $f_x(x)$  and assuming an invertible transformation  $y(x) = a + bx$  the pmf  $f_y(y)$  for the discrete scalar random variable  $y$  is given by

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

Likewise, if  $x \sim F_x$  is a discrete random vector with pmf  $f_x(x)$  and assuming an invertible transformation  $y(x) = a + Bx$  the pmf  $f_y(y)$  for the discrete random vector  $y$  is given by

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

## Density

If  $x \sim F_x$  is a continuous scalar random variable with pdf  $f_x(x)$  and assuming an invertible transformation  $y(x) = a + bx$  the pdf  $f_y(y)$  for the continuous random scalar  $y$  is given by

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

where  $|b|$  is the absolute value of  $b$ . The transformation of the corresponding differential element is

$$dy = |b| dx$$

Likewise, if  $x \sim F_x$  is a continuous random vector with pdf  $f_x(x)$  and assuming an invertible transformation  $y(x) = a + Bx$  the pdf  $f_y(y)$  for the continuous random vector  $y$  is given by

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

where  $|\det(B)|$  is the absolute value of the determinant  $\det(B)$ . The transformation of the corresponding infinitesimal volume element is

$$dy = |\det(B)| dx$$

## 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  $y \sim F_y$  is

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

and

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

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

## Importance of affine transformations

The constants  $a$  and  $B$  (or  $a$  and  $b$  in the univariate case) are the parameters of the **location-scale family**  $F_y(a, B)$  created from  $F_x$ . Many important distributions are location-scale families such as the normal distribution (cf. Section 5.3 and Section 6.3) and the location-scale *t*-distribution (Section 5.6 and Section 6.6). Furthermore, key procedures in multivariate statistics such as orthogonal transformations (including PCA) or whitening transformations (e.g. the Mahalanobis transformation) are affine transformations.

## 3.2 General invertible transformation

### Transformation rule

As above we assume  $x$  is a scalar and  $x$  is a vector and consider the general invertible transformation.

For a scalar variable the transformation is specified by  $y(x) = h(x)$  and the back-transformation by  $x(y) = h^{-1}(y)$ . For a vector this becomes  $y(x) = h(x)$  with back-transformation  $x(y) = h^{-1}(y)$ . The functions  $h(x)$  and  $h^{-1}(y)$  are assumed to be invertible.

## Probability mass function

If  $x \sim F_x$  is a discrete scalar random variable with pmf  $f_x(x)$  then the pmf  $f_y(y)$  of the transformed discrete scalar random variable  $y(x)$  is given by

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

Likewise, for a discrete random vector  $x \sim F_x$  with pmf  $f_x(x)$  the pmf  $f_y(y)$  for the discrete random vector  $y(x)$  is obtained by

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

## Density

If  $x \sim F_x$  is a continuous scalar random variable with pdf  $f_x(x)$  the pdf  $f_y(y)$  of the transformed continuous scalar random variable  $y(x)$  is given by

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

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

$$dy = |Dy(x)| dx$$

Note that  $|Dx(y)| = |Dy(x)|^{-1}|_{x=x(y)}$ .

Likewise, for a continuous random vector  $x \sim F_x$  with pdf  $f_x(x)$  the pdf  $f_y(y)$  for the continuous random vector  $y(x)$  is obtained by

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

where  $Dx(y)$  is the Jacobian matrix of the inverse transformation  $x(y)$ . The transformation of the infinitesimal volume element is

$$dy = |\det(Dy(x))| dx$$

Note that  $|\det(Dx(y))| = |\det(Dy(x))|^{-1}|_{x=x(y)}$ .

## 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  $E(x) = \mu_x$  and  $\text{Var}(x) = \Sigma_x$  are the mean and variance of the original random vector  $x$ . In the **delta method** the transformation  $y(x)$  resp.  $y(x)$  is linearised around the mean  $\mu_x$  respectively  $\mu_x$  and the mean and variance resulting from the linear transformation is reported.

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

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

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

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

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

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

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

and

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

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

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

and

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

## Invertible affine transformation as special case

The invertible affine transformation (Section 3.1) is special case of the general invertible transformation.

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 Convolution of random variables

#### Sum of independent random variables

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

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

where each  $x_i \sim F_{x_i}$  has its own distribution and corresponding pdmf  $f_{x_i}(x)$ . The corresponding means are  $E(x_i) = \mu_i$  and the variances are  $\text{Var}(x_i) = \sigma_i^2$ . As the  $x_i$  are independent, and therefore uncorrelated, the covariances  $\text{Cov}(x_i, x_j) = 0$  vanish for  $i \neq j$ .

With  $x = (x_1, \dots, x_n)^T$  and  $\mathbf{1}_n = (1, 1, \dots, 1)^T$  the relationship between  $y$  and  $x$  can be written as the linear transformation

$$y = \mathbf{1}_n^T x$$

As  $y$  is a scalar and  $x$  a vector the transformation from  $x$  to  $y$  is not invertible.

#### Moments

With  $E(x) = \mu$  and  $\text{Var}(x) = \text{Diag}(\sigma_1^2, \dots, \sigma_n^2)$  the mean of the random variable  $y$  equals

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

and the variance of  $y$  is

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

(cf. Section 3.1). Thus both the mean and variance of  $y$  are simply the sums of the individual means and variances (note that for the variance this only holds because the individual variables are uncorrelated).

#### Convolution

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

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

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

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

and (discrete case)

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

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

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

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

#### Central limit theorem

The **central limit theorem**, first postulated by [Abraham de Moivre \(1667–1754\)](#) and later proved by [Pierre-Simon Laplace \(1749–1827\)](#) asserts that the distribution of the sum of  $n$  independent and identically distributed random variables with finite mean and finite variance converges in the limit of large  $n$  to a normal distribution (Section 5.3), 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 a normal distribution.

## 4 Evaluation

### 4.1 Loss functions

#### Loss function

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

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

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

#### Risk function

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

$$R(Q, a) = \mathbb{E}_Q(L(x, a))$$

The risk is **mixture preserving in Q** meaning that

$$R(Q_\lambda, a) = (1 - \lambda)R(Q_0, a) + \lambda R(Q_1, a)$$

for the mixture  $Q_\lambda = (1 - \lambda)Q_0 + \lambda Q_1$  with  $0 < \lambda < 1$  and  $Q_0 \neq Q_1$ . This follows from the linearity of expectation.

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

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

where the expectation is replaced by the sample average.

#### Minimising risk

Minimising  $R(Q, a)$  with regard to  $a$  finds optimal predictions

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

with associated minimum risk  $R(Q, a^*)$ .

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

#### Equivalent loss functions

Adding a positive scaling factor  $k > 0$  or an additive term  $c(x)$  to a loss function generates a family of **equivalent loss functions**

$$L^{\text{equiv}}(x, a) = kL(x, a) + c(x)$$

with associated risk

$$R^{\text{equiv}}(Q, a) = kR(Q, a) + \mathbb{E}_Q(c(x))$$

Equivalent losses yield the **same risk minimiser**  $\arg \min_a R(Q, a)$  and the **same loss minimiser**  $\arg \min_a L(x, a)$  for fixed  $x$ .

## 4.2 Common loss functions

#### Squared loss

The **squared loss** or **squared error**

$$L_{\text{sq}}(x, a) = (x - a)^2$$

is one of the most commonly used loss functions. The corresponding risk is the **mean squared loss** or **mean squared error** (MSE)

$$R_{\text{sq}}(Q, a) = \mathbb{E}_Q((x - a)^2)$$

which is **minimised at the mean**  $a^* = \mathbb{E}(Q)$ . This follows from  $R_{\text{sq}}(Q, a) = \mathbb{E}_Q(x^2) - 2a \mathbb{E}_Q(x) + a^2$  and  $dR_{\text{sq}}(Q, a)/da = -2 \mathbb{E}_Q(x) + 2a$ . The **minimum risk**  $R_{\text{sq}}(a^*) = \text{Var}(Q)$  equals the **variance**.

## 0-1 loss

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

$$L_{01}(x, a) = \begin{cases} -[x = a] & \text{discrete case} \\ -\delta(x - a) & \text{continuous case} \end{cases}$$

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

$$R_{01}(Q, a) = -q(a)$$

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

## Asymmetric loss

The **asymmetric loss** can be defined as

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

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

## Absolute loss

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

$$L_{\text{abs}}(x, a) = |x - a|$$

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

## 4.3 Scoring rules

### Proper scoring rules

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

<sup>1</sup>Treating scoring rules as loss functions implies a negative orientation. However, some authors adopt the opposite convention and treat scoring rules as positively oriented utility functions.

The **risk** of  $P$  under  $Q$  is the expected score

$$R(Q, P) = E_Q(S(x, P))$$

For a **proper scoring rule**, the risk  $R(Q, P)$  is smallest when the quoted model  $P$  matches the true model  $Q$ . The minimal risk, achieved for  $P = Q$ , leads to the **properness inequality**

$$R(Q, P) \geq R(Q, Q)$$

For a **strictly proper** scoring rule, the minimum risk is realised only for the true model, so **equality holds exclusively** for  $P = Q$ .

Proper scoring rules are useful because they enable identification and approximation of data-generating distributions and their parameters via risk minimisation or minimisation of the associated scoring-rule divergences. This allows to generalise conventional statistical approaches based on the logarithmic scoring rule (Section 4.4).

### Scoring-rule entropy

The minimum risk associated with a proper scoring rule is called the **scoring-rule entropy**  $R(Q) = R(Q, Q)$ . With it the properness inequality becomes

$$R(Q, P) \geq R(Q)$$

For a **proper** scoring rule, the entropy  $R(Q)$  is **concave** in  $Q$ . For a **strictly proper** scoring rule, the entropy  $R(Q)$  is **strictly concave**. This means that

$$R(Q_\lambda) \geq (1 - \lambda)R(Q_0) + \lambda R(Q_1)$$

for the mixture  $Q_\lambda = (1 - \lambda)Q_0 + \lambda Q_1$  with  $0 < \lambda < 1$  and  $Q_0 \neq Q_1$  (for strict concavity replace  $\geq$  by  $>$ ).

This follows from the fact that the risk  $R(Q, P)$  is mixture-preserving in  $Q$ . Hence,  $R(Q_\lambda) = R(Q_\lambda, Q_\lambda) = (1 - \lambda)R(Q_0, Q_\lambda) + \lambda R(Q_1, Q_\lambda)$ . Applying properness  $R(Q_i, Q_\lambda) \geq R(Q_i)$  with  $i \in \{0, 1\}$  yields concavity.

## Scoring-rule divergence

The **scoring-rule divergence** between the distributions  $Q$  and  $P$  equals the excess risk given by

$$D(Q, P) = R(Q, P) - R(Q)$$

For a **proper** scoring rule, the divergence  $D(Q, P) \geq 0$  is always **non-negative** and with  $D(Q, P) = 0$  if  $P = Q$ . For a **strictly proper** scoring rule  $D(Q, P) = 0$  only when  $P = Q$ .

The scoring-rule divergence  $D(Q, P)$  is **convex in  $Q$**  for fixed  $P$  for a **proper** scoring rule. It is **strictly convex in  $Q$**  for a **strictly proper** scoring rule. The convexity of  $D(Q, P)$  in  $Q$  derives from the concavity of  $R(Q)$  and the fact that  $R(Q, P)$  is mixture-preserving in  $Q$ .

Divergences induced by proper scoring rules correspond to **Bregman divergences** applied to probability distributions, with the convex generator being the negative entropy.

## Equivalent scoring rules

Equivalent scoring rules

$$S^{\text{equiv}}(x, P) = kS(x, P) + c(x)$$

have associated equivalent divergences

$$D^{\text{equiv}}(Q, P) = R^{\text{equiv}}(Q, P) - R^{\text{equiv}}(Q) = kD(Q, P)$$

Thus, (strictly) proper scoring rules remain (strictly) proper under equivalence transformations. Furthermore, for  $k = 1$  equivalent scoring rules are **strongly equivalent** as their divergences are identical.

## Further properties

Proper scoring rules also enjoy several additional properties not mentioned above. For example, various **decompositions** exist for their risk, and the scoring-rule divergence satisfies a **generalised Pythagorean theorem**.

## 4.4 Common scoring rules

### Logarithmic scoring rule

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

$$S_{\log}(x, P) = -\log p(x)$$

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

$$R_{\log}(Q, P) = -E_Q \log p(x) = H(Q, P)$$

which is uniquely minimised for  $P = Q$ . Thus, the log-loss is **strictly proper**. Moreover, the log-loss is noted as the only **local** strictly proper scoring rule, as it solely depends on the value of the pdmf at the observed outcome  $x$ , and not on any other features of the distribution  $P$ .

The mean log-loss is also known as **cross-entropy** denoted by  $H(Q, P)$ .

The minimum risk (scoring-rule entropy) equals the **information entropy** denoted by  $H(Q)$ :

$$R_{\log}(Q) = -E_Q \log q(x) = H(Q)$$

The properness inequality  $H(Q, P) \geq H(Q)$ , with equality exclusively for  $P = Q$  and relating cross-entropy and information entropy, is known as **Gibbs' inequality**.

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

$$\begin{aligned} D_{\text{KL}}(Q, P) &= R_{\log}(Q, P) - R_{\log}(Q) \\ &= H(Q, P) - H(Q) \\ &= E_Q \log \left( \frac{q(x)}{p(x)} \right) \end{aligned}$$

The KL divergence is **invariant under a change of variables**. Specifically, for arbitrary one-to-one transformations between  $x$  and  $y$ , even for continuous random variables,  $D_{\text{KL}}(Q_x, P_x) = D_{\text{KL}}(Q_y, P_y)$ .

More generally, the KL divergence also satisfies the **data processing inequality** (DPI). Applying a transformation (stochastic or deterministic, possibly coarsening) that produces  $y$  from  $x$ , cannot increase the divergence, so that  $D_{\text{KL}}(Q_x, P_x) \geq D_{\text{KL}}(Q_y, P_y)$ . Note that a change of variables is a special case of data processing (with identity).

Divergences between distributions satisfying the DPI (and hence being invariant under a change of variables) form the class of *f*-divergences. The KL divergence is the *only* divergence induced by a proper scoring rule (i.e. the only Bregman divergence) that is also an *f*-divergence.

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

$$\begin{aligned} R_{\log}(\hat{Q}_n, P(\theta)) &= -\frac{1}{n} \sum_{i=1}^n \log p(x_i | \theta) \\ &= -\frac{1}{n} \ell_n(\theta) \end{aligned}$$

Minimising the empirical risk is thus equivalent to maximising the log-likelihood  $\ell_n(\theta)$ .

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

### Quadratic or Brier scoring rule

Assume the categorical distributions  $Q = \text{Cat}(q)$  with class probabilities  $q = (q_1, \dots, q_K)^T$  and  $P = \text{Cat}(p)$  with corresponding class probabilities  $p = (p_1, \dots, p_K)^T$  and 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  $x_k = 1$ .

The **quadratic scoring rule**, also known as **Brier scoring rule**, evaluates the categorical forecast  $P$  given a realisation  $x$  from the categorical distribution  $Q$  using the **squared Euclidean distance** between  $x$  and  $p$ :

$$\begin{aligned} S_{\text{Brier}}(x, P) &= \|x - p\|^2 \\ &= (x - p)^T(x - p) \\ &= 1 - 2x^T p + p^T p \\ &= 1 - 2p_k + p^T p \end{aligned}$$

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

The corresponding risk is

$$\begin{aligned} R_{\text{Brier}}(Q, P) &= E_Q(S(x, P)) \\ &= 1 - 2q^T p + p^T p \end{aligned}$$

### 4.4 Common scoring rules

which is uniquely minimised for  $P = Q$  with  $p = q$ . Thus, the Brier score is **strictly proper**.

The minimum risk (scoring-rule entropy) is

$$R_{\text{Brier}}(Q) = 1 - q^T q$$

The divergence induced by the Brier score is the **squared Euclidean distance** between  $q$  and  $p$ :

$$\begin{aligned} D_{\text{Brier}}(Q, P) &= R_{\text{Brier}}(Q, P) - R_{\text{Brier}}(Q) \\ &= (q - p)^T(q - p) \\ &= \|q - p\|^2 \end{aligned}$$

### Spherical scoring rule

The **cosine similarity** between two vectors  $a$  and  $b$  is the cosine of the angle  $\phi$  between the two vectors:

$$\text{cos\_sim}(a, b) = \cos \phi(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$$

with  $a \cdot b = a^T b$ ,  $\|a\| = (\mathbf{a}^T \mathbf{a})^{1/2}$  and  $\|b\| = (\mathbf{b}^T \mathbf{b})^{1/2}$ .

The **cosine distance** is the complement

$$\text{cos\_dist}(a, b) = 1 - \text{cos\_sim}(a, b)$$

The **spherical scoring rule** evaluates the categorical forecast  $P = \text{Cat}(p)$  given a realisation  $x$  from the categorical distribution  $Q = \text{Cat}(q)$  using the **negative cosine similarity** between the probability vectors  $x$  and  $p$ :

$$\begin{aligned} S_{\text{sph}}(x, P) &= -\text{cos\_sim}(x, p) \\ &= -x^T p / \|p\| \\ &= -p_k / \|p\| \end{aligned}$$

The spherical score is *not local* as the pmf for  $P$  is evaluated across all  $K$  classes, not just at the realised class  $k$ . Depending on  $p$ , it may range from a minimum of  $-1$  (angle  $\phi = 0$ , zero degrees) and to a maximum of  $0$  (angle,  $\phi = \pi/2$ , 90 degrees).

The corresponding risk is

$$\begin{aligned} R_{\text{sph}}(Q, P) &= E_Q(S(x, P)) \\ &= -\mathbf{q}^T \mathbf{p} / \|\mathbf{p}\| \\ &= -\|\mathbf{q}\| \cos_{\text{sim}}(\mathbf{q}, \mathbf{p}) \end{aligned}$$

which is proportional to the negative cosine similarity between  $\mathbf{q}$  and  $\mathbf{p}$ . It is uniquely minimised for  $P = Q$  with  $\mathbf{p} = \mathbf{q}$ . Therefore, the spherical scoring rule is **strictly proper**.

The minimum risk (scoring-rule entropy) is

$$R_{\text{sph}}(Q) = -\|\mathbf{q}\|$$

The divergence induced by the spherical score is

$$\begin{aligned} D_{\text{sph}}(Q, P) &= R_{\text{sph}}(Q, P) - R_{\text{sph}}(Q) \\ &= -\mathbf{q}^T \mathbf{p} / \|\mathbf{p}\| + \|\mathbf{q}\| \\ &= \|\mathbf{q}\| \cos_{\text{dist}}(\mathbf{q}, \mathbf{p}) \end{aligned}$$

and hence is proportional to the cosine distance between  $\mathbf{q}$  and  $\mathbf{p}$ .

### Proper but not strictly proper scoring rules

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

$$S_{\text{sq}}(x, P) = (x - E(P))^2$$

The corresponding risk is

$$\begin{aligned} R_{\text{sq}}(Q, P) &= E_Q((x - E(P))^2) \\ &= (E(Q) - E(P))^2 + \text{Var}(Q) \end{aligned}$$

which is minimised at  $P = Q$  but also at any distribution  $P$  with the same mean as  $Q$ .

The minimum risk (scoring-rule entropy) is the variance

$$R_{\text{sq}}(Q) = \text{Var}(Q)$$

The scoring-rule divergence is the squared distance between the two means

$$\begin{aligned} D_{\text{sq}}(Q, P) &= R_{\text{sq}}(Q, P) - R_{\text{sq}}(Q) \\ &= (E(Q) - E(P))^2 \end{aligned}$$

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

The **Dawid-Sebastiani scoring rule** is a related scoring rule given by

$$S_{\text{DS}}(x, P) = \log \text{Var}(P) + \frac{(x - E(P))^2}{\text{Var}(P)}$$

It is equivalent to the log-loss applied to a normal model  $P$ .

The corresponding risk is

$$R_{\text{DS}}(Q, P) = \log \text{Var}(P) + \frac{(E(Q) - E(P))^2}{\text{Var}(P)} + \frac{\text{Var}(Q)}{\text{Var}(P)}$$

which is minimised at  $P = Q$  but also at any distribution  $P$  with  $E(P) = E(Q)$  and  $\text{Var}(P) = \text{Var}(Q)$ .

The minimum risk (scoring-rule entropy) is

$$R_{\text{DS}}(Q) = \log \text{Var}(Q) + 1$$

The scoring-rule divergence is

$$\begin{aligned} D_{\text{DS}}(Q, P) &= R_{\text{DS}}(Q, P) - R_{\text{DS}}(Q) \\ &= \frac{(E(Q) - E(P))^2}{\text{Var}(P)} + \frac{\text{Var}(Q)}{\text{Var}(P)} - \log \left( \frac{\text{Var}(Q)}{\text{Var}(P)} \right) - 1 \end{aligned}$$

which vanishes at  $P = Q$  but also at any  $P$  for which  $E(P) = E(Q)$  and  $\text{Var}(P) = \text{Var}(Q)$ .

### Other strictly proper scoring rules

Other useful strictly proper scoring rules include:

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

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

## 5 Univariate distributions

### 5.1 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 5.1). 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

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

and the variance is

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

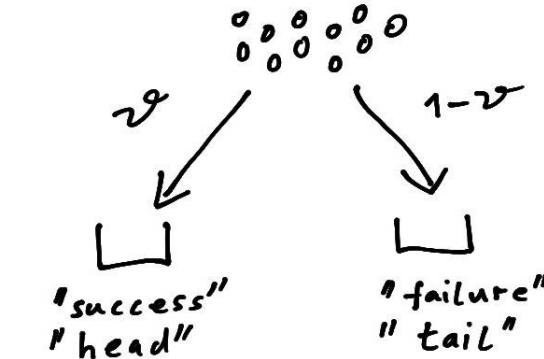


Figure 5.1: Binomial urn model.

The corresponding pmf is

$$p(x|n, \theta) = W_2 \theta^x (1 - \theta)^{n-x}$$

The binomial coefficient

$$W_2 = \binom{n}{x}$$

in the pmf accounts for the number of possible permutations of  $n$  items of two distinct types (“success” and “failure”). Note that the binomial coefficient  $W_2$  does not depend on  $\theta$ .

#### 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

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

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

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

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 3.3), for large  $n$  the binomial distribution can be well approximated by a normal distribution (Section 5.3) with the same mean and variance. This is known as the [De Moivre–Laplace theorem](#).

## 5.2 Beta distribution

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

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

### Standard parametrisation

A beta-distributed random variable is denoted by

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

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

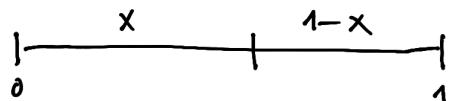


Figure 5.2: Stick breaking visualisation of a beta random variable.

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 5.2). 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

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

In this density the **beta function**

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

serves as normalisation factor.

The beta distribution can assume a number of different shapes, depending on the values of  $\alpha_1$  and  $\alpha_2$  (see Figure 5.3).

#### 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()`.

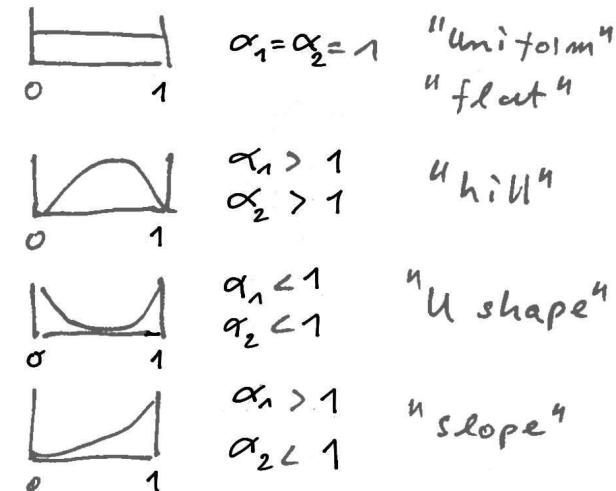


Figure 5.3: Shapes of the pdf of the beta distribution.

#### 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  $m = 2$ .

## 5.3 Normal distribution

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

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

### 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

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

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()`.

### 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 5.4 shows the pdf and cdf of the standard normal distribution.

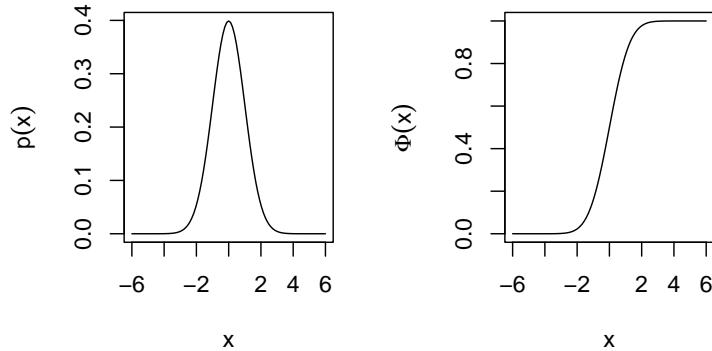


Figure 5.4: Probability density function (left) and cumulative density function (right) of the standard normal distribution.

### Special case: delta distribution

The **delta distribution**  $\delta(\mu)$  is obtained as the limit of  $N(\mu, \varepsilon\sigma^2)$  for  $\varepsilon \rightarrow 0$  and where  $\sigma^2$  is a positive number (e.g.  $\sigma^2 = 1$ ). Thus  $\delta(\mu)$  is a continuous distribution representing a point mass at  $\mu$ .

The corresponding pdf  $\delta(x|\mu)$  is called the **Dirac delta function**, even though it is not an ordinary function. It satisfies  $\delta(x|\mu) = 0$  for all  $x \neq \mu$  with an infinite spike at  $\mu$  but still integrates to one.

### 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.

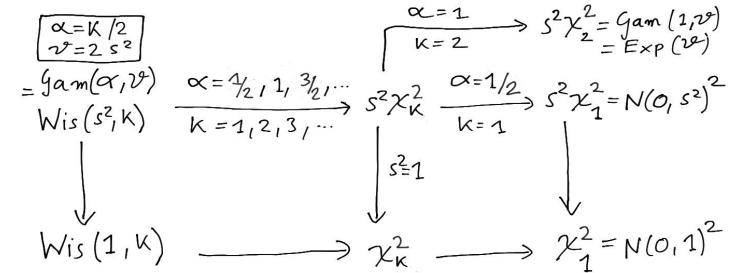


Figure 5.5: The gamma and the univariate Wishart distribution and their relatives.

### 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**.

## 5.4 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 5.5 illustrates the relationship of the gamma and the univariate Wishart distribution with these related distributions.

## Standard parametrisation

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

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

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

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

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

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

In this density the **gamma function**

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

is part of the normalisation factor.

The gamma function can also be used as continuous version of the factorial as  $\Gamma(x) = (x - 1)!$  for any positive integer  $x$ .

### R code

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

## Wishart parametrisation

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

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

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

In the above the scale parameter  $s^2$  is a 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 6.4.

In the Wishart parametrisation the mean is

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

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

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

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

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

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

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

where  $k$  is called the degree of freedom.

The chi-squared distribution has mean

$$E(x) = k$$

and variance

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

Figure 5.6 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  $z^2 \sim \chi_1^2 = \text{Wis}(1, 1) = N(0, 1)^2$ .

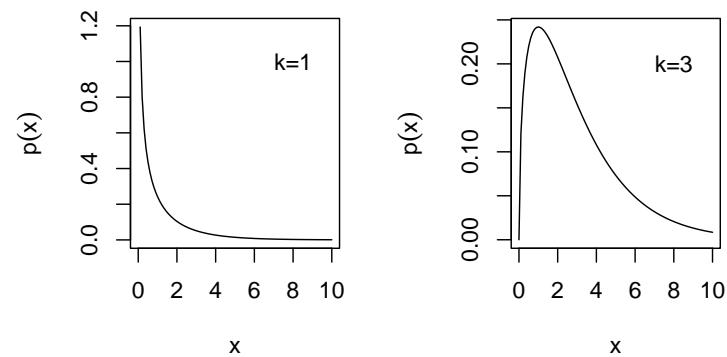


Figure 5.6: Probability density function of the chi-squared distribution.

#### R code

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

### Special case: exponential distribution

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

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

with scale parameter  $\theta$ .

It has mean

$$E(x) = \theta$$

and variance

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

and the pdf is

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

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

#### R code

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

### Scale transformation

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

Hence,

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

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

### Convolution property

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

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

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

In Wishart parametrisation this becomes

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

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

The above includes the following two specific constructions:

- If  $x_1, \dots, x_n \sim \text{Exp}(\theta)$  are independent samples from  $\text{Exp}(\theta)$  then the sum  $y = \sum_{i=1}^n x_i \sim \text{Gam}(\alpha = n, \theta)$  is gamma distributed with the same scale parameter.
- The sum of  $k$  independent scaled chi-squared random variables  $s^2 \chi_1^2$  with one degree of freedom and identical scale parameter  $s^2$  yields a scaled chi-squared random variable  $s^2 \chi_k^2$  with degree of freedom  $k$  and the same scale parameter. Thus, if  $z_1, z_2, \dots, z_k \sim N(0, 1)$  are  $k$  independent samples from  $N(0, 1)$  then  $\sum_{i=1}^k z_i^2 \sim \chi_k^2$ .

### 5.5 Inverse gamma distribution

The **inverse gamma distribution**  $\text{IGam}(\alpha, \beta)$  is a continuous distribution and is also known as **univariate inverse Wishart distribution**  $\text{IWis}(\psi, k)$  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 5.4).

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

### Standard parametrisation

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

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

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  $\alpha > 1$ )

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

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

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

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

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

#### 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()`.

#### Relation to the gamma distribution

The inverse gamma distribution is closely linked to the gamma distribution. Assume that the random variable  $y$  follows a gamma distribution with

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

then the inverse random variable  $x = 1/y$  follows an inverse gamma distribution with inverted scale parameter

$$x = \frac{1}{y} \sim \text{IGam}\left(\alpha, \beta = \frac{1}{\theta}\right)$$

where  $\alpha$  is the shared shape parameter,  $\theta$  the scale parameter of the gamma distribution and  $\beta$  the scale parameter of the inverse gamma distribution.

Correspondingly, the density  $p_x(x|\alpha, \beta)$  of the inverse gamma distribution is obtained from the density  $p_y(y|\alpha, \theta)$  of the gamma distribution via

$$p_x(x|\alpha, \beta) = \frac{1}{x^2} p_y\left(\frac{1}{x} \middle| \alpha, \theta = \frac{1}{\beta}\right)$$

#### Wishart parametrisation

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

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

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 6.5.

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

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

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

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

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

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

#### Relation to the univariate Wishart distribution

The univariate inverse Wishart distribution is closely linked to the univariate Wishart distribution. Assume that the random variable  $y$  follows an univariate Wishart distribution with

$$y \sim \text{Wis}(s^2, k)$$

then the inverse random variable  $x = 1/y$  follows a univariate inverse Wishart distribution with inverted scale parameter

$$x = \frac{1}{y} \sim \text{IWis}\left(\psi = \frac{1}{s^2}, k\right)$$

where  $k$  is the shared shape parameter,  $s^2$  the scale parameter of the univariate Wishart distribution and  $\psi$  the scale parameter of the univariate inverse Wishart distribution.

Correspondingly, the density  $p_x(x|\psi, k)$  of the univariate inverse Wishart distribution is obtained from the density  $p_y(y|s^2, k)$  of the univariate Wishart distribution via

$$p_x(x|\psi, k) = \frac{1}{x^2} p_y\left(\frac{1}{x} \middle| s^2 = \frac{1}{\psi}, k\right)$$

### 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

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

has mean (for  $\kappa > 0$ )

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

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**

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

with mean (for  $\nu > 2$ )

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

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 5.6). It is also common in Bayesian analysis.

### Special case: inverse chi-squared distribution

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

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

where  $k$  is called the degree of freedom.

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

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

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

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

### Relation to the chi-squared distribution

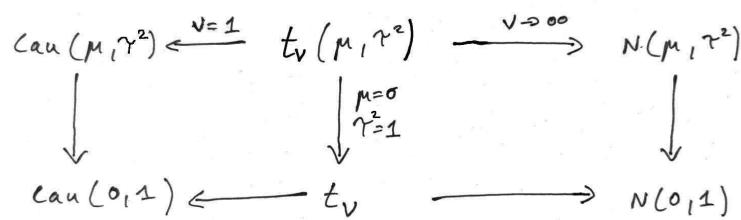
The inverse chi-squared distribution is closely linked to the chi-squared distribution. Assume that the random variable  $y$  follows a chi-squared distribution with

$$y \sim \chi_k^2$$

then the inverse random variable  $x = 1/y$  follows an inverse chi-squared distribution

$$x = \frac{1}{y} \sim \chi_k^{-2}$$

where  $k$  is the shared degree of freedom (shape parameter).

Figure 5.7: The location-scale  $t$ -distribution and its relatives.

### Scale transformation

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

Hence,

- $\beta \text{IGam}(\alpha, 1) = \text{IGam}(\alpha, \beta)$ ,
- $\psi \text{IWis}(1, k) = \text{IWis}(\psi, k)$ ,
- $\kappa\mu \text{IWis}(1, k = \kappa + 2) = \text{IWis}(\psi = \kappa\mu, k = \kappa + 2)$  and
- $v\tau^2 \text{IWis}(1, k = v) = \text{IWis}(\psi = v\tau^2, k = v)$

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

The scaled inverse chi-squared distribution is frequently used in the biased mean parametrisation with  $\tau = \psi/v$  and  $v = k$ . Then  $\psi \chi_k^{-2}$  is equal to  $v\tau^2 \chi_v^{-2} = \text{IWis}(\psi = v\tau^2, k = v)$  which is sometimes also written as  $\chi^{-2}(v, \tau^2)$ .

## 5.6 Location-scale $t$ -distribution

The **location-scale  $t$ -distribution**  $t_v(\mu, \tau^2)$  is a continuous distribution and is a generalisation of the normal distribution  $N(\mu, \tau^2)$  (Section 5.3) with an additional parameter  $v > 0$  (degrees of freedom) controlling the probability mass in the tails.

Special cases include the **Student's  $t$ -distribution**  $t_v$ , the **normal distribution**  $N(\mu, \tau^2)$  and the **Cauchy distribution**  $\text{Cau}(\mu, \tau^2)$ . Figure 5.7

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

### Standard parametrisation

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

$$x \sim t_v(\mu, \tau^2)$$

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

The mean is (for  $v > 1$ )

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

and the variance (for  $v > 2$ )

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

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

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

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 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**  $t_v = t_v(0, 1)$ . 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  $v > 1$ ) and variance  $\text{Var}(x) = \frac{v}{v-2}$  (for  $v > 2$ ).

The pdf of  $t_v$  is

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

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  $v \rightarrow \infty$  the location-scale  $t$ -distribution  $t_v(\mu, \tau^2)$  reduces to the **normal distribution**  $N(\mu, \tau^2)$  (Section 5.3). Correspondingly, for  $v \rightarrow \infty$  the Student's  $t$ -distribution becomes equal to the **standard normal distribution**  $N(0, 1)$ .

See Section 6.6 for further details.

### Special case: Cauchy distribution

For  $v = 1$  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  $v = 1$  (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  $x \sim t_v(\mu, \tau^2)$  then  $y = w(x - \mu) \sim t_v$ . This location-scale transformation reduces a location-scale  $t$ -distributed random variable to a Student's  $t$ -distributed random variable.

Conversely, if  $y \sim t_v$  then  $x = \mu + \tau y \sim t_v(\mu, \tau^2)$ . 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 5.3, 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

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

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

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

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.

## 6 Multivariate distributions

### 6.1 Multinomial distribution

The **multinomial distribution**  $\text{Mult}(n, \theta)$  is the multivariate generalisation of the binomial distribution  $\text{Bin}(n, \theta)$  (Section 5.1) from two classes to  $K$  classes.

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

#### Standard parametrisation

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

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

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

The multinomial distribution is best illustrated by an urn model distributing  $n$  items into  $K$  bins where  $\theta$  contains the corresponding bin probabilities (Figure 6.1).

The expected value is

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

The covariance matrix is

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

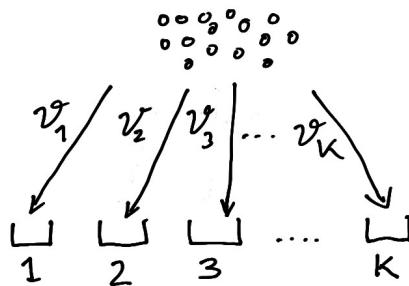


Figure 6.1: Multinomial urn model.

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

The corresponding pmf is

$$p(x|\theta) = W_K \prod_{k=1}^K \theta_k^{x_k}$$

The multinomial coefficient

$$W_K = \binom{n}{x_1, \dots, x_K}$$

in the pmf accounts for the number of possible permutations of  $n$  items of  $K$  distinct types. Note that the multinomial coefficient  $W_K$  does not depend on  $\theta$ .

While all  $K$  elements of  $x$  appear in the pmf recall that due the dependencies among the  $x_k$  the pmf is defined over a  $K - 1$  dimensional support.

For  $K = 2$  the multinomial distribution reduces to the binomial distribution (Section 5.1).

#### 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  $\mu_k \in [0, n]$  and  $\mu^T \mathbf{1}_K = \sum_{k=1}^K \mu_k = n$ , 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

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

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  $x_k \in \{0, 1\}$  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 indicating the class to which the item has been allocated. This is called “one hot encoding”, as opposed to “integer encoding” (i.e. stating the class number directly).

The expected value is

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

The covariance matrix is

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

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

The corresponding pmf is

$$p(x|\theta) = \prod_{k=1}^K \theta_k^{x_k} = \begin{cases} \theta_k & \text{if } x_k = 1 \\ 0 & \text{otherwise} \end{cases}$$

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

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

## Convolution property

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

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

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

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

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

## 6.2 Dirichlet distribution

The **Dirichlet distribution**  $\text{Dir}(\alpha)$  is the multivariate generalisation of the beta distribution  $\text{Beta}(\alpha_1, \alpha_2)$  (Section 5.2) 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

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

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



Figure 6.2: 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 6.2). Thus, the  $x_k$  may be used as the exclusive proportions or probabilities for  $K$  classes.

The mean is

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

and the variance is

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

The covariance matrix is singular by construction because of the dependencies among the elements of  $x$ . 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  $\text{Dir}(\alpha)$  is

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

In this density the normalisation factor is given by the **multivariate beta function**

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

For  $K = 2$  it reduces to the conventional beta function  $B(\alpha_1, \alpha_2)$ .

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

For  $K = 2$  the Dirichlet distribution reduces to the beta distribution (Section 5.2).

#### 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  $\mu_k \in [0, 1]$  and  $\mu^T \mathbf{1}_K = \sum_{k=1}^K \mu_k = 1$ , 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

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

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  $x$ . In component notation it is

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

### Special case: symmetric Dirichlet distribution

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

### Special case: uniform distribution

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

## 6.3 Multivariate normal distribution

The **multivariate normal distribution**  $N(\mu, \Sigma)$  generalises the univariate normal distribution  $N(\mu, \sigma^2)$  (Section 5.3) from one to  $d$  dimensions.

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

### Standard parametrisation

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

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

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

with mean

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

and variance

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

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

The pdf is given by

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

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

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

For  $d = 1$  the random vector  $\mathbf{x} = x$  is a scalar,  $\boldsymbol{\mu} = \mu$ ,  $\boldsymbol{\Sigma} = \sigma^2$  and thus the multivariate normal distribution reduces to the univariate normal distribution (Section 5.3).

### 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  $\boldsymbol{\Sigma}$  it is less obvious how to define a suitable matrix scale parameter.

Let  $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$  be the eigendecomposition of the positive definite matrix  $\boldsymbol{\Sigma}$ . Then  $\boldsymbol{\Sigma}^{1/2} = \mathbf{U}\boldsymbol{\Lambda}^{1/2}\mathbf{U}^T$  is the principal matrix square root and  $\boldsymbol{\Sigma}^{-1/2} = \mathbf{U}\boldsymbol{\Lambda}^{-1/2}\mathbf{U}^T$  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  $\mathbf{W} = \mathbf{Q}\boldsymbol{\Sigma}^{-1/2}$  is called a **whitening matrix** based on  $\boldsymbol{\Sigma}$  and  $\mathbf{L} = \mathbf{W}^{-1} = \boldsymbol{\Sigma}^{1/2}\mathbf{Q}^T$  is the corresponding **inverse whitening matrix**. By construction, the matrix  $\mathbf{L}$  provides a factorisation of the covariance matrix by  $\mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}$ . Similarly,  $\mathbf{W}$  factorises the precision matrix by  $\mathbf{W}^T \mathbf{W} = \boldsymbol{\Sigma}^{-1}$ . 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  $\boldsymbol{\Sigma}$  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  $\boldsymbol{\Sigma}^{-1}$  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  $|\det(\mathbf{W})| = \det(\boldsymbol{\Sigma})^{-1/2}$  and  $|\det(\mathbf{L})| = \det(\boldsymbol{\Sigma})^{1/2}$  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  $\boldsymbol{\Sigma}$  alone but there is rotational freedom due to  $\mathbf{Q}$ . Standard choices are

- $\mathbf{Q}^{\text{ZCA}} = \mathbf{I}$  for ZCA-type factorisation with  $\mathbf{W}^{\text{ZCA}} = \boldsymbol{\Sigma}^{-1/2}$  and
- $\mathbf{Q}^{\text{PCA}} = \mathbf{U}^T$  for PCA-type factorisation with  $\mathbf{W}^{\text{PCA}} = \boldsymbol{\Lambda}^{-1/2}\mathbf{U}^T$ . 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  $\boldsymbol{\Sigma}$ , 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  $\mathbf{Q}^{\text{Chol}} = (\mathbf{L}^{\text{Chol}})^T \boldsymbol{\Sigma}^{-1/2}$ .

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  $\boldsymbol{\Sigma} = \mathbf{V}^{1/2} \mathbf{P} \mathbf{V}^{1/2}$ ) 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**  $N(0, I)$  has mean  $\mu = 0$  and variance  $\Sigma = I$ . The corresponding pdf is

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

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

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

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

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

### Special case: multivariate delta distribution

The **multivariate delta distribution**  $\delta(\mu)$  is obtained as the limit of  $N(\mu, \varepsilon A)$  for  $\varepsilon \rightarrow 0$  and where  $A$  is a positive definite matrix (e.g.  $A = I$ ). Thus  $\delta(\mu)$  is a continuous distribution representing a point mass at  $\mu$ .

The corresponding pdf  $\delta(x|\mu)$  is called the **multivariate Dirac delta function**, even though it is not an ordinary function. It satisfies  $\delta(x|\mu) = 0$  for all  $x \neq \mu$  with an infinite spike at  $\mu$  but still integrates to one.

### Location-scale transformation

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

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

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

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

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

### Convolution property

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

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

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

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

## 6.4 Wishart distribution

The Wishart distribution  $\text{Wis}(S, k)$  is a multivariate generalisation of the gamma distribution  $\text{Gam}(\alpha, \theta)$  (Section 5.4) 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 6.5), otherwise  $X$  is singular and positive semi-definite.

The distribution has mean

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

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

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

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

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

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  $\mathbf{X}$ .

Part of the normalisation factor in the density is the **multivariate gamma function**

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

For  $d = 1$  it reduces to the standard gamma function. The multivariate gamma function also occurs in densities of other matrix-variate distributions.

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 5.4).

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  $\mathbf{z}\mathbf{z}^T \sim \text{Wis}(\mathbf{S}, 1)$ .



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

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

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  $\mathbf{S} = \mathbf{I}$  the Wishart distribution reduces to the **standard Wishart distribution**

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

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

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

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

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

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

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

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

The **Bartlett decomposition** of the standard multivariate Wishart  $\text{Wis}(\mathbf{I}, k)$  distribution for any real  $k > d - 1$  is obtained by Cholesky factorisation of the random matrix  $\mathbf{X} = \mathbf{Z}\mathbf{Z}^T$ . By construction  $\mathbf{Z}$  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  $\mathbf{Z}$  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  $\mathbf{X} = \mathbf{Z}\mathbf{Z}^T \sim \text{Wis}(\mathbf{I}, k)$ .

## 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  $Y = WXW^T \sim \text{Wis}(I, k)$  as  $WSW^T = I$ . 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  $X = LYL^T \sim \text{Wis}(S, k)$  as  $LIL^T = S$ . 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 to the partial restriction of the shape parameter  $k$  to integer values the multivariate Wishart distribution is **not infinitely divisible** for  $d > 1$ .

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

## 6.5 Inverse Wishart distribution

The **inverse Wishart distribution**  $\text{IWis}(\Psi, k)$  is a multivariate generalisation of the inverse gamma distribution  $\text{IGam}(\alpha, \beta)$  (Section 5.5) from one to  $d$  dimensions. It is linked to the Wishart distribution  $\text{Wis}(S, k)$  (Section 6.4).

### Standard parametrisation

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

$$X \sim \text{IWis}(\Psi, k)$$

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  $\text{IWis}(\Psi, k)$  has pdf

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

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

If  $\Psi$  is a scalar  $\psi$  (and  $d = 1$ ) then the multivariate inverse Wishart distribution reduces to the univariate inverse Wishart distribution (Section 5.5).

### R code

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

### Relation to the Wishart distribution

The inverse Wishart distribution is closely linked to the Wishart distribution. Assume that the random variable  $Y$  (positive definite and symmetric) follows a Wishart distribution with

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

then the inverse random variable  $X = Y^{-1}$  (positive definite and symmetric) follows an inverse Wishart distribution with inverted scale parameter

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

where  $k$  is the shared shape parameter,  $S$  the scale parameter of the Wishart distribution and  $\Psi$  the scale parameter of the inverse Wishart distribution.

Correspondingly, the density  $p_X(X|\Psi, k)$  of the inverse Wishart distribution is obtained from the density  $p_Y(Y|S, k)$  of the Wishart distribution via

$$p_X(X|\Psi, k) = \det(X)^{-d-1} p_Y(X^{-1}|S = \Psi^{-1}, k)$$

The exponent  $-d - 1$  in the Jacobian determinant for matrix inversion arises because  $X$  and  $Y$  are symmetric (for non-symmetric matrices the exponent is  $-2d$ ).

### 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

$$X \sim \text{IWis}(\Psi = \kappa M, k = \kappa + d + 1)$$

with mean (for  $\kappa > 0$ )

$$\mathbb{E}(X) = M$$

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**

$$X \sim \text{IWis}(\Psi = \nu T, k = \nu + d - 1)$$

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

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

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  $X \sim \text{IWis}(\Psi, k)$  then the scaled symmetric random matrix  $AXA^T$  is also inverse Wishart distributed with  $AXA^T \sim \text{IWis}(A\Psi A^T, k)$  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  $\Psi$ .

## 6.6 Multivariate t-distribution

The **multivariate t-distribution**  $t_v(\mu, T)$  is a multivariate generalisation of the location-scale t-distribution  $t_v(\mu, \tau^2)$  (Section 5.6) from one to  $d$  dimensions. It is a generalisation of the multivariate normal distribution  $N(\mu, T)$  (Section 6.3) with an additional parameter  $v > 0$  (degrees of freedom) controlling the probability mass in the tails.

Special cases include the **multivariate standard t-distribution**  $t_v(0, I)$ , 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

$$x \sim t_v(\mu, T)$$

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  $v > 0$  prescribes the degrees of freedom. For small values of  $v$  the distribution is heavy-tailed and as a result only moments of order smaller than  $v$  are finite and defined.

The mean is (for  $v > 1$ )

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

and the variance (for  $v > 2$ )

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

The pdf of  $t_v(\mu, T)$  is

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

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

For  $d = 1$  the random vector  $x = x$  is a scalar,  $\mu = \mu$ ,  $T = \tau^2$  and thus the multivariate t-distribution reduces to the location-scale t-distribution (Section 5.6).

### R code

The `mnormt` package implements the multivariate t-distribution. The function `mnormt::dmt()` provides the pdf and `mnormt::pmnt()` 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  $L$  in place of a matrix dispersion parameter  $T$ .

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

Note that  $T$  alone does not fully determine  $L$  and  $W$  due to rotational freedom, see the discussion in Section 6.3 for details.

### Special case: multivariate standard t-distribution

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

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

The pdf of  $t_v(0, I)$  is

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

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

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

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 6.3). Correspondingly, for  $\nu \rightarrow \infty$  the multivariate standard  $t$ -distribution  $t_\nu(0, I)$  becomes equal to the **multivariate standard normal distribution**  $N(0, I)$ .

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

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

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

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

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

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

### Special case: multivariate Cauchy distribution

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

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

It has pdf

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

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

### Special case: multivariate standard Cauchy distribution

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

It has pdf

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

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

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

Conversely, if  $y \sim t_\nu(0, I)$  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  $x = \mu + Ly$  with  $\text{Var}(y) = \nu/(\nu - 2)I$  we get  $\text{Cov}(x, y) = \nu/(\nu - 2)L$ . This provides a means to choose between different factorisations of  $T$  and  $T^{-1}$ . For example, if positive correlation between corresponding elements in  $x$  and  $y$  is desired then the diagonal elements in  $L$  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  $x \sim Cau(\mu, T)$  then  $y = W(x - \mu) \sim Cau(0, I)$ . Conversely, if  $y \sim Cau(0, I)$  then  $x = \mu + Ly \sim Cau(\mu, T)$ .

## 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 6.3, and the multivariate Cauchy distribution ( $\nu = 1$ ) with the additional restriction that the dispersion parameters are proportional.

For the Cauchy distribution with  $T_i = a_i^2 T$ , where  $a_i > 0$  are positive scalars,

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

## 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

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

and let  $x|z$  be multivariate normal

$$x|z \sim N(\mu, \Sigma = zT)$$

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

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

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

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

and let

$$x|Z \sim N(\mu, \Sigma = Z)$$

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

# 7 Exponential families

## 7.1 Definition of an exponential family

### Exponential tilting

A distribution family  $P(\eta)$  for a random variable  $x$  is an **exponential family** if it is generated by **exponential tilting** of a base distribution  $B$ , resulting in a pdmf of the form

$$p(x|\eta) = \underbrace{h(x)}_{\text{base}} \underbrace{e^{\langle \eta, t(x) \rangle}}_{\text{exponential tilt}} / \underbrace{z(\eta)}_{\text{normaliser}} \\ = h(x) e^{\langle \eta, t(x) \rangle - a(\eta)}$$

where

- $t(x)$  are the **canonical statistics**,
- $\eta$  are the **canonical parameters**,
- $h(x)$  is a positive **base function** (typically unnormalised),
- $z(\eta)$  is the **partition function** and
- $a(\eta) = \log z(\eta)$  the corresponding **log-partition function**.

The base pdmf is obtained at  $\eta = 0$  yielding  $b(x) = p(x|\eta = 0) = h(x)/z(0)$ . If  $h(x)$  is already a normalised pdmf then  $z(0) = 1$  and  $b(x) = h(x)$ .

The above presentation of exponential families assumes a univariate random variable (scalar  $x$ ) but also applies to multivariate random variables (vector  $x$  or matrix  $X$ ).

Likewise, canonical statistics and parameters are written as vectors but these may also be scalars or matrices (or a combination of both). The use of inner product notation  $\langle , \rangle$  includes all these cases, recall that for scalars  $\langle a, b \rangle = ab$ , for vectors  $\langle a, b \rangle = a^T b$  and for matrices  $\langle A, B \rangle = \text{Tr}(A^T B) = \text{Vec}(A)^T \text{Vec}(B)$ .

## Canonical statistics

The canonical statistics  $t(x)$  are transformations of  $x$ , usually simple functions such as the identity ( $x$ ), the square ( $x^2$ ), the inverse ( $1/x$ ) or the logarithm ( $\log x$ ).

Typically, the number of canonical statistics and hence the dimension of  $t$  is small.

The canonical statistics  $t(x)$  may be affinely dependent. If this is the case there is a vector  $\eta_0$  for which

$$\langle \eta_0, t(x) \rangle = \text{const.}$$

A common example is when  $x$  is a vector of counts  $(n_1, \dots, n_K)^T$  for  $K$  classes with a fixed total count  $n = \sum_{i=1}^K n_k = x^T \mathbf{1}_K$  and the canonical statistics are  $t(x) = x$  and thus include all  $K$  counts.

If the elements in  $t(x)$  are affinely independent the representation of the exponential family is **minimal** or **complete**, otherwise the representation is **non-minimal** or **overcomplete**.

## Canonical parameters and identifiability

For each canonical statistic there is a corresponding canonical parameter so the dimensions and shape of  $t(x)$  and  $\eta$  match.

In a minimal representation the canonical parameters of the exponential family are **identifiable** and hence distinct parameter settings for  $\eta$  yield distinct distributions.

Conversely, in a non-minimal or overcomplete representation there are redundant elements in the canonical parameters  $\eta$  and the distributions within the exponential family are **not identifiable**. Specifically, there will be multiple  $\eta$  yielding the same underlying distribution.

In the example above, where  $x$  is a vector of counts with a fixed total count and  $t(x) = x$  and corresponding canonical parameters  $\eta$ , the effective number of parameter is  $K - 1$  rather than  $K$ , so there is one redundant parameter in  $\eta$ .

## Moment and cumulant generating functions

The moment generating function for the canonical statistics  $t(x)$  is

$$\begin{aligned} M(\tau) &= E\left(e^{\langle \tau, t(x) \rangle}\right) \\ &= \int_x e^{\langle \tau, t(x) \rangle} p(x|\eta) dx \\ &= \int_x e^{\langle \tau, t(x) \rangle} e^{\langle \eta, t(x) \rangle} h(x)/z(\eta) dx \\ &= \left( \int_x e^{\langle \tau + \eta, t(x) \rangle} h(x) dx \right) / z(\eta) \\ &= z(\tau + \eta) / z(\eta) \end{aligned}$$

Correspondingly, the cumulant generating function is

$$\begin{aligned} K(\tau) &= \log M(\tau) \\ &= a(\tau + \eta) - a(\eta) \end{aligned}$$

Thus, the moment and cumulant generating functions for the canonical statistics are closely linked to the partition and log-partition functions, respectively.

## 7.2 Roles of the partition function

### Normalising factor

The pdmf  $p(x|\eta)$  must integrate to one. Therefore, given  $h(x)$  and  $t(x)$  the partition function  $z(\eta)$  is obtained by

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

For discrete  $x$  replace the integral by a sum.

Consequently partition function  $z(\eta)$  is also called the **normaliser** and the log-partition function  $a(\eta)$  is called the **log-normaliser**.

Constructed as weighted integral of exponentials of linear functions, the partition function  $z(\eta)$  and the log-partition function  $a(\eta)$  are **convex** with regard to  $\eta$ .

For an exponential family in **minimal representation** the (log)-partition function(s) are **strictly convex**.

## Definition of parameter space

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

## Moments of canonical statistics

The first cumulant (the mean) and second cumulant (the variance) are obtained as the first and second derivatives of the cumulant generating function  $K(\tau) = a(\tau + \eta) - a(\eta)$  evaluated at  $\tau = 0$ , respectively. As a result, the log-partition function  $a(\eta)$  provides a practical way to obtain the mean and variance of the canonical statistics  $t(x)$ .

Specifically, computing its gradient yields the mean

$$\begin{aligned} E(t(x)) &= \mu_t = \nabla a(\eta) \\ &= \frac{\nabla z(\eta)}{z(\eta)} \end{aligned}$$

and computing the Hessian matrix the covariance matrix

$$\begin{aligned} \text{Var}(t(x)) &= \Sigma_t = \nabla \nabla^T a(\eta) \\ &= \frac{\nabla \nabla^T z(\eta)}{z(\eta)} - \left( \frac{\nabla z(\eta)}{z(\eta)} \right) \left( \frac{\nabla z(\eta)}{z(\eta)} \right)^T \end{aligned}$$

For an exponential family with minimal representation the log-partition function is strictly convex. Furthermore, the variance  $\Sigma_t$  is a positive definite matrix and invertible.

For overcomplete representations the log-partition function is convex but not strictly convex. Furthermore, the covariance matrix is positive semi-definite and not invertible.

By construction, the log-partition function  $a(\eta)$  is finite in the interior of its parameter space. Therefore *all* moments and cumulants of the canonical statistics  $t(x)$  exist (are finite) and are given by the derivatives of  $a(\eta)$  at  $\eta$ .

## 7.3 Further properties

### Equivalent representations

An exponential family admits many equivalent representations such that different specifications of canonical statistics  $t(x)$  and base function  $h(x)$  describe the same family.

First, any invertible linear transformation of the canonical statistic  $t(x)$  yields the same distribution family.

Second, any member of the family, say  $P(\eta_0)$ , can serve as its base distribution. Specifically, with  $p(x|\eta_0)$  used as base the pdmf for the exponential family  $P(\eta)$  is

$$p(x|\eta) = e^{\langle \eta - \eta_0, t(x) \rangle - (a(\eta) - a(\eta_0))} p(x|\eta_0)$$

which is in exponential family form.

Third, the base function  $h(x)$  can be left unnormalised, so there are infinitely many positive base functions  $h(x)$  that yield the same base pdmf  $b(x)$  after normalisation.

Fourth, any factors in  $h(x)$  of the form  $e^{\langle \eta_0, t(x) \rangle}$  for some constant  $\eta_0$  can be removed from  $h(x)$  and absorbed into the exponential by replacing  $\eta$  with  $\eta + \eta_0$  (thus leading to a different set of canonical parameters).

As a result, for many (but not all) commonly used exponential families the base function can be set to  $h(x) = 1$  or some other constant value, so that all dependence on  $x$  enters through the canonical statistics  $t(x)$  via  $\langle \eta, t(x) \rangle$ .

### Natural exponential families

If the canonical statistic is the identity,  $t(x) = x$  or  $t(x) = x$ , the family is a **natural exponential family (NEF)**. A univariate NEF has a scalar canonical statistic  $t(x)$  and a scalar canonical parameter  $\eta$  and is thus a one-parameter family.

## Alternative parametrisations

An exponential family can be parametrised by three different sets of parameters:

- 1) **canonical parameters**  $\eta$ ,
- 2) **expectation parameters**  $\mu_t = E(t(x))$  (the mean of the canonical statistics  $t(x)$ ), as well as
- 3) **conventional parameters**  $\theta$  (such as mean and variance of  $x$ ).

If the exponential family is minimal then there is a one-to-one map between the canonical parameters  $\eta$  and the expectation parameters  $\mu_t$ .

The canonical and the expectation parameters can be expressed as a function of the conventional parameters  $\theta$ .

Often, some expectation parameters  $\mu_t$  correspond to conventional parameters  $\theta$  (e.g., if one of the canonical statistics is  $x$ , then the corresponding expectation parameter is the mean of  $x$ ).

## Exponential families and change of variables

Assume that  $x$  is a random variable with an exponential-family distribution and  $y(x)$  is an invertible transformation (see Section 3.2).

The resulting pdmf for  $y$  is

$$p(y|\eta) = h_y(y) e^{\langle \eta, t_y(y) \rangle} / z(\eta)$$

with transformed canonical statistics

$$t_y(y) = t_x(x(y))$$

For a discrete random variable the base function changes to

$$h_y(y) = h_x(x(y))$$

whereas for a continuous random variable the base function becomes

$$h_y(y) = |Dx(y)| h_x(x(y))$$

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

Thus, for both discrete and continuous random variables the exponential-family form is preserved under a change of variables.

## 7.4 Univariate exponential families

Table 7.1 lists univariate exponential families, more details about these distributions are found in Chapter 5.

$\text{Bin}(n, \theta)$  is a one-parameter exponential family with  $n$  assumed fixed.

For  $\text{Bin}(n, \theta)$  and  $\text{Ber}(\theta)$  the conventional parameter is

$$\theta = \text{logit}^{-1}(\eta) = \frac{e^\eta}{1 + e^\eta}$$

(logistic function). Conversely, since this is a one-to-one map, the canonical parameter equals

$$\eta = \text{logit}(\theta) = \log\left(\frac{\theta}{1 - \theta}\right)$$

Apart from  $\text{Bin}(n, \theta)$  all families listed in Table 7.1 have constant base function ( $h(x) = 1$ ).

Furthermore,  $\text{Bin}(n, \theta)$ ,  $\text{Ber}(\theta)$  and  $\text{Exp}(\theta)$  are NEFs.  $N(\mu, \sigma^2)$  with fixed  $\sigma^2$  (variance),  $\text{Gam}(\alpha, \theta)$  with fixed  $\alpha$  (shape) and  $\text{Wis}(s^2, k)$  with fixed  $k$  (shape) are also NEFs.

Table 7.1: Common univariate exponential families

Distribution	$h(x)$	$z(\eta)$	$\eta$	$t(x)$	$\mu_t$
Bin( $n, \theta$ )	$W_2$	$(1 + e^\eta)^n$	$\text{logit}(\theta)$	$x$	$\theta$
Ber( $\theta$ )	1	$1 + e^\eta$	$\text{logit}(\theta)$	$x$	$\theta$
Beta( $\alpha_1, \alpha_2$ )	1	$B(\eta_1 + 1, \eta_2 + 1)$	$\begin{pmatrix} \alpha_1 - 1 \\ \alpha_2 - 1 \end{pmatrix}$	$\begin{pmatrix} \log x \\ \log(1-x) \end{pmatrix}$	$\begin{pmatrix} \psi^{(0)}(\alpha_1) - \psi^{(0)}(m) \\ \psi^{(0)}(\alpha_2) - \psi^{(0)}(m) \end{pmatrix}$
$N(\mu, \sigma^2)$	1	$(-\pi\eta_2^{-1})^{1/2}$ $\exp(-\frac{1}{4}\eta_1^2\eta_2^{-1})$	$\begin{pmatrix} \sigma^{-2}\mu \\ -\frac{1}{2}\sigma^{-2} \end{pmatrix}$	$\begin{pmatrix} x \\ x^2 \end{pmatrix}$	$\begin{pmatrix} \mu \\ \sigma^2 + \mu^2 \end{pmatrix}$
Gam( $\alpha, \theta$ )	1	$(-\eta_1)^{-\eta_2-1}$ $\Gamma(\eta_2 + 1)$	$\begin{pmatrix} -1/\theta \\ \alpha - 1 \end{pmatrix}$	$\begin{pmatrix} x \\ \log x \end{pmatrix}$	$\begin{pmatrix} \alpha\theta \\ \psi^{(0)}(\alpha) + \log \theta \end{pmatrix}$
Exp( $\theta$ )	1	$-\eta^{-1}$	$-1/\theta$	$x$	$\theta$
Wis( $s^2, k$ )	1	$(-\eta_1)^{-\eta_2-1}$ $\Gamma(\eta_2 + 1)$	$\begin{pmatrix} -\frac{1}{2}s^{-2} \\ \frac{k}{2} - 1 \end{pmatrix}$	$\begin{pmatrix} x \\ \log x \end{pmatrix}$	$\begin{pmatrix} ks^2 \\ \psi^{(0)}(\frac{k}{2}) + \log(2s^2) \end{pmatrix}$
IGam( $\alpha, \beta$ )	1	$(-\eta_1)^{\eta_2+1}$ $\Gamma(-\eta_2 - 1)$	$\begin{pmatrix} -\beta \\ -\alpha - 1 \end{pmatrix}$	$\begin{pmatrix} x^{-1} \\ \log x \end{pmatrix}$	$\begin{pmatrix} \alpha/\beta \\ -\psi^{(0)}(\alpha) + \log \beta \end{pmatrix}$
IWis( $\psi, k$ )	1	$(-\eta_1)^{\eta_2+1}$ $\Gamma(-\eta_2 - 1)$	$\begin{pmatrix} -\frac{\psi}{2} \\ -\frac{k}{2} - 1 \end{pmatrix}$	$\begin{pmatrix} x^{-1} \\ \log x \end{pmatrix}$	$\begin{pmatrix} k/\psi \\ -\psi^{(0)}(\frac{k}{2}) + \log(\frac{\psi}{2}) \end{pmatrix}$

Notes:

- $W_2 = \binom{n}{x}$  is the binomial coefficient.
- $B(\alpha_1, \alpha_2)$  is the beta function.
- $m = \alpha_1 + \alpha_2$ .
- $\psi^{(0)}(x) = \frac{d}{dx} \log \Gamma(x)$  is the **digamma function**.

## 7.5 Multivariate exponential families

Table 7.2 lists multivariate exponential families, more details about these distributions are found in Chapter 6.

Mult( $n, \theta$ ) is a  $(K - 1)$ -parameter exponential family with  $n$  assumed fixed.

For Mult( $n, \theta$ ) and Cat( $\theta$ ) the conventional parameters are given by

$$\boldsymbol{\theta} = \text{softmax}(\boldsymbol{\eta}) = \frac{(\exp \eta_k)}{\sum_{i=1}^K \exp \eta_i}$$

As the softmax function is invariant against translation and is a many-to-one map, its inverse is not unique and the canonical parameters

$$\boldsymbol{\eta} = (c + \log \theta_k)$$

are determined by the conventional parameters  $\boldsymbol{\theta}$  only up to a constant  $c$ . This representation using  $K$  canonical parameters  $\boldsymbol{\eta}$  is non-minimal, hence  $\boldsymbol{\eta}$  is not identifiable and different values of  $\boldsymbol{\eta}$  can represent the same distribution.

A minimal representation with  $K - 1$  parameters  $\eta_1, \dots, \eta_{K-1}$  and  $\eta_K = 0$  corresponds to  $c = -\log \theta_K$  and  $\eta_k = \log(\theta_k / \theta_K)$ . For  $K = 2$  this yields the minimal representations of Bin( $n, \theta$ ) and Ber( $\theta$ ) shown in Table 7.1.

Apart from Mult( $n, \theta$ ) all families listed in Table 7.2 have constant base function ( $h(x) = 1$ ).

Furthermore, Mult( $n, \theta$ ) and Cat( $\theta$ ) are NEFs.  $N(\mu, \Sigma)$  with fixed  $\Sigma$  (variance) and Wis( $S, k$ ) with fixed  $k$  (shape) are NEFs as well.

Table 7.2: Common multivariate exponential families

Distribution	$h(x)$	$z(\eta)$	$\eta$	$t(x)$	$\mu_t$
Mult( $n, \theta$ )	$W_K$	$(\sum_{k=1}^K \exp \eta_k)^n$	$(c + \log \theta_k)$	$x$	$\theta$
Cat( $\theta$ )	1	$\sum_{k=1}^K \exp \eta_k$	$(c + \log \theta_k)$	$x$	$\theta$
Dir( $\alpha$ )	1	$B(\eta + 1)$	$(\alpha_k - 1)$	$(\log x_k)$	$(\psi^{(0)}(\alpha_k) - \psi^{(0)}(m))$
$N(\mu, \Sigma)$	1	$\det(-\pi \eta_2^{-1})^{1/2} \exp(-\frac{1}{4} \eta_1^T \eta_2^{-1} \eta_1)$	$\begin{pmatrix} \Sigma^{-1} \mu \\ -\frac{1}{2} \Sigma^{-1} \end{pmatrix}$	$\begin{pmatrix} x \\ xx^T \end{pmatrix}$	$\begin{pmatrix} \mu \\ \Sigma + \mu \mu^T \end{pmatrix}$
Wis( $S, k$ )	1	$\frac{\det(-\eta_1)^{-\eta_2 - \frac{d+1}{2}}}{\Gamma_d(\eta_2 + \frac{d+1}{2})}$	$\begin{pmatrix} -\frac{1}{2} S^{-1} \\ \frac{k}{2} - \frac{d+1}{2} \end{pmatrix}$	$\begin{pmatrix} X \\ \log \det(X) \end{pmatrix}$	$\begin{pmatrix} kS \\ \psi_d^{(0)}(\frac{k}{2}) + \log \det(2S) \end{pmatrix}$
IWis( $\Psi, k$ )	1	$\frac{\det(-\eta_1)^{\eta_2 + \frac{d+1}{2}}}{\Gamma_d(-\eta_2 - \frac{d+1}{2})}$	$\begin{pmatrix} -\frac{1}{2} \Psi \\ -\frac{k}{2} - \frac{d+1}{2} \end{pmatrix}$	$\begin{pmatrix} X^{-1} \\ \log \det(X) \end{pmatrix}$	$\begin{pmatrix} k\Psi^{-1} \\ -\psi_d^{(0)}(\frac{k}{2}) + \log \det(\frac{\Psi}{2}) \end{pmatrix}$

Notes:

- $W_K = \binom{n}{x_1, \dots, x_K}$  is the multinomial coefficient for  $K$  groups.
- $B(\alpha)$  is the multivariate beta function.
- $m = \sum_{i=k}^K \alpha_k$ .
- $\psi_d^{(0)} = \frac{d}{dx} \log \Gamma_d(x) = \sum_{i=1}^d \psi^{(0)}(x - (i - 1)/2)$  is the **multivariate digamma function**.

See also: [Exponential family \(Wikipedia\)](#) and Efron (2022).

## Bibliography

- Efron, B. 2022. *Exponential Families in Theory and Practise*. Cambridge University Press. <https://doi.org/10.1017/9781108773157>.
- 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>.