Probability and Distribution Refresher

Korbinian Strimmer

25 March 2024

Table of contents

| We | Upd Lice | e 1 ates 1 nse 1 | |
|----|---|--|--|
| Pr | e face Abo Abo | ut the author 2 ut the notes 2 | |
| 1 | Com 1.1 1.2 1.3 1.4 | binatorics3Some basic mathematical notation3Number of permutations3De Moivre-Sterling approximation of the factorial4Multinomial and binomial coefficient4 | |
| 2 | Prot 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 2.10 | ability6Random variables6Probability mass and density function7Distribution function and quantile function7Families of distributions8Expectation of a random variable9Jensen's inequality for the expectation9Probability as expectation9Moments and variance of a random variable10Random vectors and their mean and variance10Correlation matrix11 | |
| 3 | Tran 3.1 3.2 3.3 3.4 | sformations12Affine or location-scale transformation of random variables12General invertible transformation of random variables13Exponential tilting and exponential families15Sums of random variables and convolution16 | |
| 4 | Univ 4.1 4.2 4.3 4.4 4.5 | ariate distributions17Bernoulli distribution17Binomial distribution17Beta distribution19Normal distribution21Gamma distribution and special cases22 | |

Table of contents

| | 4.6 | Inverse gamma distribution | 26 | | |
|----|--------------|---|----|--|--|
| | 4.7 | Location-scale <i>t</i> -distribution and special cases | 28 | | |
| 5 | Mult | tivariate distributions | 32 | | |
| | 5.1 | Categorical distribution | 32 | | |
| | 5.2 | Multinomial distribution | 33 | | |
| | 5.3 | Dirichlet distribution | 34 | | |
| | 5.4 | Multivariate normal distribution | 36 | | |
| | 5.5 | Wishart distribution | 37 | | |
| | 5.6 | Inverse Wishart distribution | 38 | | |
| Bi | Bibliography | | | | |

Welcome

The Probability and Distribution Refresher notes were written by Korbinian Strimmer from 2018–2024. This version is from 25 March 2024.

If you have any questions, comments, or corrections please get in touch!¹

Updates

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

• online version of the Probability and Distribution Refresher notes.

You may also wish to download the Probability and Distribution Refresher notes as

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

License

These notes are licensed to you under Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

Preface

About the author

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

About the notes

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

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

This includes the currently offered modules:

- MATH27720 Statistics 2: Likelihood and Bayes and
- MATH38161 Multivariate Statistics

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

• MATH20802 Statistical Methods.

¹Email address: korbinian.strimmer@manchester.ac.uk

The factorial can also be obtained using the gamma function

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

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

1.3 De Moivre-Sterling approximation of the factorial

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

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

or equivalently on logarithmic scale

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

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

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

1.4 Multinomial and binomial coefficient

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

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

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

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

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

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

1 Combinatorics

1.1 Some basic mathematical notation

Summation:

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

Multiplication:

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

Indicator function:

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

Scalar: plain type, typically lower case (x, θ) , sometimes upper case (K).

Vector: bold type, lower case (x, θ).

Matrix: bold type, upper case (X, Σ).

1.2 Number of permutations

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

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

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

0! = 1

as there is exactly one permutation of zero objects.

1 Combinatorics

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

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

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

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

2 Probability

2.1 Random variables

A **random variable** describes a random experiment. The set of all possible outcomes is the **sample space** or **state space** of the random variable and is denoted by $\Omega = \{\omega_1, \omega_2, ...\}$. The outcomes ω_i are the **elementary events**. The sample space Ω can be finite or infinite. Depending on type of outcomes the random variable is **discrete** or **continuous**.

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

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

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

This implies

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

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

From the above it is clear that the definition and theory of probability is closely linked to set theory, and in particular to measure theory. Indeed, viewing probability as a special type of measure allows for an elegant treatment of both discrete and continuous random variables. 2 Probability

2.2 Probability mass and density function

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

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

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

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

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

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

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

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

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

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

2.3 Distribution function and quantile function

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

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



Figure 2.1: Density function and distribution function.

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

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

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

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

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

2.4 Families of distributions

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

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

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

If for any pair of different parameter values $\theta_1 \neq \theta_2$ we get distinct distributions with $F_{\theta_1} \neq F_{\theta_2}$ then the distribution family F_{θ} is said to be **identifiable** by the parameter θ .

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

2.5 Expectation of a random variable

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

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

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

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

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

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

2.6 Jensen's inequality for the expectation

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

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

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

2.7 Probability as expectation

Probability itself can also be understood as an expectation. For an event *A* we can define a corresponding indicator function $1_{x \in A}$ for an elementary element *x* to be part of *A*. From the above it then follows

$$\mathrm{E}(1_{x\in A})=\mathrm{Pr}(A)\,,$$

Interestingly, one can develop the whole theory of probability from this perspective (e.g., Whittle 2000).

2.8 Moments and variance of a random variable

The moments of a random variable are defined as follows:

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

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

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

A distribution does not necessarily need to have any finite first or higher moments. An example is the location-scale *t*-distribution (Section 4.7) that depending on the value of the parameter ν may not have a mean or variance (or other higher moments).

2.9 Random vectors and their mean and variance

In addition to scalar random variables we often make use of random vectors and also random matrices. $^{\rm 2}$

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

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

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

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

2 Probability

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

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

2.10 Correlation matrix

The **correlation matrix** *P* ("upper case rho", not "upper case p") is the variance standardised version of the covariance matrix Σ .

Specifically, denote by V the diagonal matrix containing the variances

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

then the correlation matrix *P* is given by

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

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

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

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

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

3 Transformations

3.1 Affine or location-scale transformation of random variables

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

y = a + bx

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

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

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

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

$$y = a + Bx$$

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

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

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

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

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

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

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

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

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

and variance

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

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

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

$$\mathbf{E}(\boldsymbol{y}) = \boldsymbol{a} + \boldsymbol{B}\,\boldsymbol{\mu}_{\boldsymbol{x}}$$

and

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

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

The constants *a* and *B* (or *a* and *b* in the univariate case) are the parameters of the **location-scale family** F_y created from F_x . Many important distributions are location-scale families such as the normal distribution (cf. Section 5.4 and Section 5.4) and the location-scale *t*-distribution (Section 4.7). Furthermore, key procedures in multivariate statistics such as orthogonal transformations (including PCA) or whitening transformations (e.g. the Mahalanobis transformation) are affine transformations.

3.2 General invertible transformation of random variables

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

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

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

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

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

3 Transformations

Likewise, for a continuous random vector x with density $f_x(x)$ the density for y is obtained by

$$f_{\boldsymbol{y}}(\boldsymbol{y}) = |\det\left(D\boldsymbol{x}(\boldsymbol{y})\right)| f_{\boldsymbol{x}}\left(\boldsymbol{x}(\boldsymbol{y})\right)$$

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

The mean and variance of the transformed random variable can typically only be approximated. Assume that $E(x) = \mu_x$ and $Var(x) = \sigma_x^2$ are the mean and variance of the original random variable *x* and $E(x) = \mu_x$ and $Var(x) = \Sigma_x$ 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 μ_x respectively μ_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 μ_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 μ_x .

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

$$\mathrm{E}(y)\approx y\left(\mu_x\right)$$

and

$$\operatorname{Var}(y) \approx \left(Dy\left(\mu_{x}\right)\right)^{2} \sigma_{x}^{2}$$

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

$$\mathrm{E}(\boldsymbol{y}) \approx \boldsymbol{y} \left(\boldsymbol{\mu}_{\boldsymbol{x}} \right)$$

and

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

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

3 Transformations

3.3 Exponential tilting and exponential families

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

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

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

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

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

The distribution family P_{η} obtained by exponential tiling is called an **exponential family**. The corresponding log-pdmf is

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

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

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

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

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

3.4 Sums of random variables and convolution

Suppose we have a sum of n independent and identically distributed (iid) random variables.

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

where each $x_i \sim F_x$ with density or probability mass function $f_x(x)$. The density or probability mass function for y is obtained by repeated application of **convolution** (symbolised by the * operator):

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

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

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

and (discrete case)

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

Convolution is commutative and associative so it can be applied in any order to compute the convolution of multiple functions. Furthermore, the convolution of probability densities / mass function yields another probability density / mass function.

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

See also: list of convolutions of probability distributions.

The **central limit theorem**, first postulated by Abraham de Moivre (1667–1754) and later proved by Pierre-Simon Laplace (1749–1827) asserts that, under appropriate conditions, the distribution of the sum of independent and identically distributed random variables converges in the limit of large n to a normal distribution (Section 4.4), even if the individual random variables are not normal. In other words, it asserts that for large n the convolution of n identical distributions typically converges to the normal distribution.

4 Univariate distributions

4.1 Bernoulli distribution

The **Bernoulli distribution** $Ber(\theta)$ is the simplest of all distribution families. It is named after Jacob Bernoulli (1655-1705) who also discovered the law of large numbers.

It describes a discrete binary random variable with two states x = 0 ("failure") and x = 1 ("success"), where the parameter $\theta \in [0, 1]$ is the probability of "success". Often the Bernoulli distribution is also referred to as "coin tossing" model with the two outcomes "heads" and "tails".

Correspondingly, the probability mass function of $Ber(\theta)$ is

$$p(x = 0|\theta) = \Pr(\text{"failure"}|\theta) = 1 - \theta$$

and

$$p(x = 1|\theta) = \Pr("success"|\theta) = \theta$$

A compact way to write the pmf of the Bernoulli distribution is

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

The log-pmf is

 $\log p(x|\theta) = x \log \theta + (1-x) \log(1-\theta)$

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

$$x \sim \operatorname{Ber}(\theta)$$
.

The expected value is $E(x) = \theta$ and the variance is $Var(x) = \theta(1 - \theta)$.

4.2 Binomial distribution

Closely related to the Bernoulli distribution is the **binomial distribution** $Bin(n, \theta)$ which results from repeating a Bernoulli experiment *n* times and counting the number of successes among the *n* trials (without keeping track of the ordering of the experiments). Thus, if x_1, \ldots, x_n are *n* independent $Ber(\theta)$ random variables then $y = \sum_{i=1}^{n} x_i$ is distributed as $Bin(n, \theta)$.



Figure 4.1: Binomial urn model.

If a random variable *y* follows the binomial distribution we write

 $y \sim \operatorname{Bin}(n, \theta)$

The corresponding probability mass function is:

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

with support $y \in \{0, 1, 2, ..., n\}$. The binomial coefficient $\binom{n}{y}$ is needed to account for the multiplicity of ways (orderings of samples) in which we can observe y successes.

The expected value is $E(y) = n\theta$ and the variance is $Var(y) = n\theta(1 - \theta)$.

If we standardise the support of the binomial variable to the unit interval with $\frac{y}{n} \in \{0, \frac{1}{n}, ..., 1\}$ then the mean is $\mathbb{E}\left(\frac{y}{n}\right) = \theta$ and the variance is $\operatorname{Var}\left(\frac{y}{n}\right) = \frac{\theta(1-\theta)}{n}$.

For n = 1 the binomial distribution reduces to the Bernoulli distribution (Section 4.1).

The binomial distribution may be illustrated by an urn model distributing n items into two bins (Figure 4.1).

As a result of the central limit theorem, the binomial distribution, obtained as the convolution of *n* Bernoulli distributions, can for large *n* be well approximated by a normal distribution (this is known as the De Moivre–Laplace theorem).

🥊 R code

The probability mass function of the binomial distribution is given by dbinom(), the cumulative distribution function is pbinom() and the quantile function is qbinom(). The binomial coefficient is computed by choose().

4.3 Beta distribution

Standard parameterisation

A beta-distributed random variable is denoted by

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

where the support is $x \in [0, 1]$ and $\alpha > 0$ and $\beta > 0$ are two shape parameters.

The density of the beta distribution $\text{Beta}(\alpha, \beta)$ is

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

This depends on the beta function defined as

$$B(z_1,z_1)=\frac{\Gamma(z_1)\Gamma(z_2)}{\Gamma(z_1+z_2)}$$

The beta distribution is very flexible and can assume a number of different shapes, depending on the value of α and β . For example, for $\alpha = \beta = 1$ it becomes the uniform distribution over the unit interval (see Figure 4.2).

A beta random variable can be visualised as breaking a unit stick of length one into two pieces of length x and 1 - x (Figure 4.3).

💡 R code

The probability density function of the beta distribution is given by dbeta(), the cumulative distribution function is pbeta() and the quantile function is qbeta().



Figure 4.2: Shapes of the density of the beta distribution.



Figure 4.3: Stick breaking visualisation of a beta random variable.

Mean parametrisation

and

Instead of employing α and β as parameters another useful reparametrisation Beta(μ , k) of the beta distribution is in terms of a mean parameter $\mu \in [0, 1]$ and a concentration parameter k > 0. These are given by

$$k = \alpha + \beta$$

 $\mu = \frac{\alpha}{\alpha + \beta}$

The original parameters can be recovered by $\alpha = \mu k$ and $\beta = (1 - \mu)k$.

The mean and variance of the beta distribution expressed in terms of μ and k are

 $E(x) = \mu$

and

$$\operatorname{Var}(x) = \frac{\mu(1-\mu)}{k+1}$$

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

The uniform distribution (with $\alpha = \beta = 1$) corresponds to $\mu = 1/2$ and k = 2.

Finally, note that the mean and variance of the continuous beta distribution closely match those of the unit-standardised discrete binomial distribution above.

4.4 Normal distribution

The **normal distribution** is the most important continuous probability distribution. It is also called **Gaussian distribution** named after Carl Friedrich Gauss (1777–1855).

The univariate normal distribution $N(\mu, \sigma^2)$ has two parameters μ (location) and σ^2 (scale) and support $x \in] -\infty, \infty[$.

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

with mean

and variance

$$Var(x) = \sigma^2$$

 $E(x) = \mu$

Probability density function (pdf):

$$p(x|\mu,\sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

The standard normal distribution is N(0, 1) with mean 0 and variance 1. The cumulative distribution function (cdf) of the standard normal N(0, 1) is

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

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

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



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

💡 R code

The normal probability density function is given by dnorm(), the cumulative distribution function is pnorm() and the quantile function is qnorm().

4.5 Gamma distribution and special cases

The gamma distribution is widely used in statistics, and also appears in various parametrisations and under some other names, such as univariate Wishart and scaled chi-squared distribution

Standard parametrisation

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

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

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

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

and variance

 $Var(x) = \alpha \theta^2$

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 probability density function (pdf) is:

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

💡 R code

The density of the gamma distribution is available in the function dgamma(). The cumulative density function is pgamma() and the quantile function is qgamma().

Wishart parametrisation and scaled chi-squared distribution

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

 $W_1\left(s^2,k\right)$

named after John Wishart (1898–1954). In the Wishart parametrisation the mean is

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

and the variance

$$Var(x) = 2ks^4$$

Another name for the one-dimensional Wishart distribution with exactly the same parametrisation is **scaled chi-squared distribution** denoted as

 $s^2\chi_k^2$

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

$$W_1\left(s^2=\frac{\mu}{k},k\right)$$

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

 $E(x) = \mu$

and the variance

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

Construction as sum of squared normals

x

A gamma distributed variable with k = 1, 2, 3... or equivalently $\alpha = 1/2, 1, 3/2, ...$ can be constructed as follows. Assume *k* independent normal random variables with mean 0 and variance s^2 :

$$z_1, z_2, \ldots, z_k \sim N(0, s^2)$$

Then the sum of the squares

$$x = \sum_{i=1}^{k} z_i^2$$

follows the distribution

$$\sim s^2 \chi_k^2$$

= W₁ (s², k)
= Gam ($\alpha = \frac{k}{2}, \theta = 2s^2$)

Chi-squared distribution

The **chi-squared distribution** χ_k^2 is a special one-parameter restriction of the gamma resp. Wishart distribution obtained when setting $s^2 = 1$ or, equivalently, $\theta = 2$ or $\mu = k$.

It has mean E(x) = k and variance Var(x) = 2k. The chi-squared distribution χ_k^2 equals $Gam(\alpha = k/2, \theta = 2)$ and $W_1(1, k)$.

Figure 4.5 shows plots the density of the chi-squared distribution for degrees of freedom k = 1 and k = 3.

💡 R code

The density of the chi-squared distribution is given by dchisq(). The cumulative density function is pchisq() and the quantile function is qchisq().



Figure 4.5: Density of the chi-squared distribution.

Exponential distribution

The **exponential distribution** $Exp(\theta)$ with scale parameter θ is another special one-parameter restriction of the gamma distribution with shape parameter set to $\alpha = 1$ (or equivalently k = 2).

The exponential distribution $\text{Exp}(\theta)$ equals $\text{Gam}(\alpha = 1, \theta)$ and $W_1(s^2 = \theta/2, k = 2)$.

The density of the exponential distribution is

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

with mean $E(x) = \theta$ and variance $Var(x) = \theta^2$.

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

💡 R code

The command dexp() returns the density of the exponential distribution, pexp() is the corresponding cumulative density function and qexp() is the quantile function.

4.6 Inverse gamma distribution

Also know as inverse univariate Wishart distribution.

Standard parametrisation

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

$$x \sim \text{Inv-Gam}(\alpha, \beta)$$

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

The inverse of x is then gamma distributed

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

where α is the shared shape parameter and θ the scale parameter of the gamma distribution.

The inverse gamma distribution Inv-Gam(α , β) has density

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

The mean of the inverse gamma distribution is

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

and the variance

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

Thus, for the mean to exist we have the restriction $\alpha > 1$ and for the variance to exist $\alpha > 2$.

💡 R code

The extraDistr package implements the inverse gamma distribution. The function extraDistr::dinvgamma() provides the density, the function extraDistr::pinvgamma() returns the corresponding cumulative density function and extraDistr::qinvgamma() is the quantile function.

4 Univariate distributions

Wishart parametrisation

The inverse gamma distribution is frequently used with a different set of parameters $\psi = 2\beta$ (scale parameter) and $\nu = 2\alpha$ (shape parameter), or conversely $\alpha = \nu/2$ and $\beta = \psi/2$. In this form it is called **one-dimensional inverse Wishart distribution**

 $W_1^{-1}(\psi, \nu)$

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

with mean given by

for $\nu > 2$ and variance

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

for $\nu > 4$.

The inverse univariate Wishart and univariate Wishart distributions are linked. If a random variable x is inverse Wishart distributed

$$x \sim W_1^{-1}(\psi, \nu)$$

then the inverse of *x* is Wishart distributed with inverted scale parameter:

$$\frac{1}{x} \sim W_1(s^2 = \psi^{-1}, k = \nu)$$

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

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

$$W_1^{-1}(\psi = \kappa \mu, \nu = \kappa + 2)$$

has mean

$$E(x) = \mu$$

with $\kappa > 0$ and the variance is

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

with $\kappa > 2$. This **mean parametrisation** is useful when employing the inverse gamma distribution as prior and posterior.

Finally, with $W_1^{-1}(\psi = \nu \tau^2, \nu)$, where $\tau^2 = \mu \frac{\kappa}{\kappa+2} = \frac{\psi}{\nu}$ is a biased mean parameter, we get the **scaled inverse chi-squared distribution** $\tau^2 \text{Inv-} \chi_{\nu}^2$ with

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

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

for $\nu > 2$ and

for $\nu > 4$.



Figure 4.6: The location-scale *t* distribution and its relatives.

4.7 Location-scale *t*-distribution and special cases

Location-scale *t*-distribution

The location-scale *t*-distribution $lst(\mu, \tau^2, \nu)$ is a generalisation of the normal distribution. It has an additional parameter $\nu > 0$ (degrees of freedom) that controls the probability mass in the tails. For small values of ν the distribution is heavy-tailed — indeed so heavy that for $\nu \leq 1$ even the mean is not defined and for $\nu \leq 2$ the variance is undefined.

The probability density of $lst(\mu, \tau^2, \nu)$ is

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

with support $x \in] - \infty, \infty[$. The mean is (for $\nu > 1$)

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

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

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

...

For $\nu \to \infty$ the location-scale *t*-distribution $lst(\mu, \tau^2, \nu)$ becomes the normal distribution $N(\mu, \tau^2)$.

Figure 4.6 illustrates the relationship of the location-scale *t* distribution $lst(\mu, \tau^2, \nu)$ with related distributions such as the normal distribution $N(\mu, \tau^2)$, Student's *t*-distribution t_{ν} and the Cauchy distribution Cau (μ, τ) discussed further below.

💡 R code

The package extraDistr implements the location-scale *t*-distribution. The function extraDistr::dlst() returns the density, extraDistr::plst() is the corresponding cumulative density function and extraDistr::qlst() is the quantile function.

Location-scale *t*-distribution as compound distribution

Suppose that

$$x|s^2 \sim N(\mu,s^2)$$

with corresponding density $p(x|s^2)$ and mean $E(x|s^2) = \mu$ and variance $Var(x|s^2) = s^2$.

Now let the variance s^2 be distributed as univariate inverse gamma / inverse Wishart

$$S^2 \sim W_1^{-1}(\psi = \kappa \sigma^2, \nu = \kappa + 2) = W_1^{-1}(\psi = \tau^2 \nu, \nu)$$

with corresponding density $p(s^2)$ and mean $E(s^2) = \sigma^2 = \tau^2 \nu / (\nu - 2)$. Note we use here both the mean parametrisation (σ^2 , κ) and the inverse chi-squared parametrisation (τ^2 , ν).

The joint density for x and s^2 is $p(x, s^2) = p(x|s^2)p(s^2)$. We are interested in the marginal density for x:

$$p(x) = \int p(x, s^2) ds^2 = \int p(s^2) p(x|s^2) ds^2$$

This is a compound distribution of a normal with fixed mean μ and variance s^2 varying according the inverse gamma distribution. Calculating the integral results in the location-scale *t*-distribution with parameters

$$x \sim \operatorname{lst}\left(\mu, \sigma^2 \frac{\kappa}{\kappa+2}, \kappa+2\right) = \operatorname{lst}\left(\mu, \tau^2, \nu\right)$$

 $E(x) = \mu$

with mean

and variance

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

From the law of total expectation and variance we can also directly verify that

$$\mathbf{E}(x) = \mathbf{E}(\mathbf{E}(x|s^2)) = \mu$$

and

$$Var(x) = E(Var(x|s^{2})) + Var(E(x|s^{2}))$$
$$= E(s^{2}) = \sigma^{2}$$
$$= \tau^{2} \frac{\nu}{\nu - 2}$$

Student's *t*-distribution

For $\mu = 0$ and $\tau^2 = 1$ the location-scale *t*-distribution becomes the Student's *t*-distribution t_{ν} . It is named after "Student" which was the pseudonym of William Sealy Gosset (1876–1937).

It has mean 0 (for $\nu > 1$) and variance $\frac{\nu}{\nu-2}$ (for $\nu > 2$).

It can thus be viewed as a generalisation of the standard normal distribution N(0, 1).

If $y \sim t_v$ then $x = \mu + \tau y$ is distributed as $x \sim \operatorname{lst}(\mu, \tau^2, v)$.

For $\nu \to \infty$ the *t*-distribution becomes equal to N(0, 1).

The probability density of t_{ν} is

$$p(x|\nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu}\,\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-(\nu+1)/2} \label{eq:px}$$

with support $x \in] - \infty, \infty[$.

💡 R code

The command dt() returns the density of the *t*-distribution, pt() is the corresponding cumulative density function and qt() is the quantile function.

Cauchy and standard Cauchy distribution

For v = 1 the location-scale *t*-distribution becomes the Cauchy distribution Cau(μ , τ) with density $p(x|\mu, \tau) = \frac{\tau}{\pi(\tau^2 + (x-\mu)^2)}$. It is named after Augustin-Louis Cauchy (1789–1857).

For $\nu = 1$ the *t*-distribution becomes the standard Cauchy distribution Cau(0, 1) with density $p(x) = \frac{1}{\pi(1+x^2)}$.

🥊 R code

The command dcauchy() returns the density of the Cauchy distribution, pcauchy() is the corresponding cumulative density function and qcauchy() is the quantile function.

5 Multivariate distributions

5.1 Categorical distribution

The **categorical distribution** is a generalisation of the Bernoulli distribution from two classes to *K* classes.

The categorical distribution $Cat(\pi)$ describes a discrete random variable with K states ("categories", "classes", "bins") where the parameter vector $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)^T$ specifies the probability of each of class so that $Pr(\text{"class } k") = \pi_k$. The parameters satisfy $\pi_k \in [0, 1]$ and $\sum_{k=1}^K \pi_k = 1$, hence there are K - 1 independent parameters in a categorical distribution (and not K).

There are two main ways to numerically represent "class k":

- i) by "integer encoding", i.e. by the corresponding integer *k*.
- ii) by "one hot encoding", i.e. by an indicator vector $\mathbf{x} = (x_1, \dots, x_K)^T = (0, 0, \dots, 1, \dots, 0)^T$ containing zeros everywhere except for the element $x_k = 1$ at position *k*. Thus all $x_k \in \{0, 1\}$ and $\sum_{k=1}^K x_k = 1$.

In the following we use "one hot encoding". Therefore sampling from a categorical distribution with parameters π

 $x \sim \operatorname{Cat}(\pi)$

yields a random index vector *x*.

The corresponding probability mass function (pmf) can be written conveniently in terms of x_k as

$$p(\boldsymbol{x}|\boldsymbol{\pi}) = \prod_{k=1}^{K} \pi_k^{x_k} = \left\{ \pi_k \quad \text{if } x_k = 1 \right\}$$

and the log-pmf as

$$\log p(\boldsymbol{x}|\boldsymbol{\pi}) = \sum_{k=1}^{K} x_k \log \pi_k = \left\{ \log \pi_k \quad \text{if } x_k = 1 \right.$$

In order to be more explicit that the categorical distribution has K - 1 and not K parameters we rewrite the log-density with $\pi_K = 1 - \sum_{k=1}^{K-1} \pi_k$ and $x_K = 1 - \sum_{k=1}^{K-1} x_k$

$$\log p(\mathbf{x}|\pi) = \sum_{k=1}^{K-1} x_k \log \pi_k + x_K \log \pi_K$$
$$= \sum_{k=1}^{K-1} x_k \log \pi_k + \left(1 - \sum_{k=1}^{K-1} x_k\right) \log \left(1 - \sum_{k=1}^{K-1} \pi_k\right)$$

Note that there is no particular reason to choose π_K as dependent of the probabilities of the other classes, in its place any other of the π_k may be selected.

The expected value is $E(x) = \pi$, in component notation $E(x_k) = \pi_k$. The covariance matrix is $Var(x) = Diag(\pi) - \pi \pi^T$, which in component notation is $Var(x_i) = \pi_i(1 - \pi_i)$ and $Cov(x_i, x_j) = -\pi_i \pi_j$.

The form of the categorical covariance matrix follows directly from the definition of the variance $Var(x) = E(xx^T) - E(x)E(x)^T$ and noting that $x_i^2 = x_i$ and $x_ix_j = 0$ if $i \neq j$. Furthermore, the categorical covariance matrix is singular by construction, as the *K* random variables x_1, \ldots, x_K are dependent through the constraint $\sum_{k=1}^{K} x_k = 1$.

For K = 2 the categorical distribution reduces to the Bernoulli Ber(θ) distribution, with $\pi_1 = \theta$ and $\pi_2 = 1 - \theta$ (Section 4.1).

5.2 Multinomial distribution

The **multinomial distribution** Mult(n, π) arises from repeated categorical sampling, in the same fashion as the binomial distribution arises from repeated Bernoulli sampling. Thus, if $x_1, ..., x_n$ are n independent Cat(π) random categorical variables then $y = \sum_{i=1}^{n} x_i$ is distributed as Mult(n, π).

The corresponding pmf describes the probability of a pattern y_1, \ldots, y_K of samples distributed across *K* classes (with $n = \sum_{k=1}^{K} y_k$):

$$p(\boldsymbol{y}|n,\theta) = \binom{n}{y_1,\ldots,y_n} \prod_{k=1}^K \pi_k^{y_k}$$

where $\binom{n}{v_1,\ldots,v_n}$ is the multinomial coefficient.

The expected value is

 $E(y) = n\pi$

which in component notation is $E(y_k) = n\pi_k$. The covariance matrix is

$$\operatorname{Var}(\boldsymbol{y}) = n(\operatorname{Diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^T)$$

which in component notation is $Var(x_i) = n\pi_i(1 - \pi_i)$ and $Cov(x_i, x_j) = -n\pi_i\pi_j$.



Figure 5.1: Multinomial urn model.

Standardised to unit interval we get:

$$\frac{y_i}{n} \in \left\{0, \frac{1}{n}, \frac{2}{n}, \dots, 1\right\}$$
$$E\left(\frac{y}{n}\right) = \pi$$
$$Var\left(\frac{y}{n}\right) = \frac{Diag(\pi) - \pi\pi^T}{n}$$
$$Var\left(\frac{y_i}{n}\right) = \frac{\pi_i(1 - \pi_i)}{n}$$
$$Cov\left(\frac{y_i}{n}, \frac{y_j}{n}\right) = -\frac{\pi_i\pi_j}{n}$$

For n = 1 the multinomial distribution reduces to the categorical distribution (Section 5.1).

For K = 2 the multinomial distribution reduces to the Binomial distribution (Section 4.2).

The multinomial distribution may be illustrated by an urn model distributing n balls into K bins (Figure 5.1).

5.3 Dirichlet distribution

Standard parametrisation

The Dirichlet distribution is the multivariate generalisation of the beta distribution. It is named after Peter Gustav Lejeune Dirichlet (1805–1859).

A Dirichlet distributed random vector is denoted by

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

with parameter $\boldsymbol{\alpha} = (\alpha_1, ..., \alpha_K)^T > 0$ and $K \ge 2$ and where the support of \boldsymbol{x} is the K - 1 dimensional simplex with $x_i \in [0, 1]$ and $\sum_{i=1}^K x_i = 1$.

The density of the Dirichlet distribution $Dir(\alpha)$ is

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

This depends on the beta function with multivariate argument defined as

$$B(z) = \frac{\prod_{k=1}^{K} \Gamma(z_k)}{\Gamma\left(\sum_{k=1}^{K} z_k\right)}$$



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

A Dirichlet random variable can be visualised as breaking a unit stick into *K* individual pieces of lengths x_1 to x_K adding up to one (Figure 5.2).

For K = 2 the Dirichlet distribution reduces to the beta distribution (Section 4.3).

Mean parametrisation

Instead of employing *a* as parameter vector another useful reparametrisation $\text{Dir}(\pi, k)$ of the Dirichlet distribution is in terms of a mean parameter π , with $\pi_i \in [0, 1]$ and $\sum_{i=1}^{K} \pi_i = 1$, and a concentration parameter k > 0. These are given by

$$k = \sum_{i=1}^{K} \alpha_i$$

 $\pi = \frac{\alpha}{k}$

and



The mean and variance of the Dirichlet distribution expressed in terms of π and k are $E(x) = \pi$

and

$$\operatorname{Var}(x) = \frac{\operatorname{Diag}(\pi) - \pi \pi^{T}}{k+1}$$

which in component notation is

$$\operatorname{Var}(x_i) = \frac{\pi_i(1 - \pi_i)}{k + 1}$$

and

$$\operatorname{Cov}(x_i, x_j) = -\frac{\pi_i \pi_j}{k+1}$$

Finally, note that the mean and variance of the continuous Dirichlet distribution closely match those of the unit-standardised discrete multinomial distribution above.

5.4 Multivariate normal distribution

The univariate normal distribution for a random scalar *x* generalises to the **multi-variate normal distribution** for a random vector $\mathbf{x} = (x_1, x_2, ..., x_d)^T$.

If *x* follows a multivariate normal distribution we write

$$\boldsymbol{x} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where μ is the mean (location) parameter and Σ the variance (scale) parameter. The corresponding density is

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \det(2\pi\boldsymbol{\Sigma})^{-1/2}\exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{T}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

As $\det(2\pi\Sigma)^{-1/2} = \det(2\pi I_d)^{-1/2} \det(\Sigma)^{-1/2} = (2\pi)^{-d/2} \det(\Sigma)^{-1/2}$ the density can also be written as

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

i.e. with explicit occurrence of the dimension *d*.

The expectation of *x* is $E(x) = \mu$ and the variance is $Var(x) = \Sigma$.

For d = 1 the random vector $\mathbf{x} = \mathbf{x}$ is a scalar and $\boldsymbol{\mu} = \boldsymbol{\mu}$ and $\boldsymbol{\Sigma} = \sigma^2$ and the multivariate normal density reduces to the univariate normal density (Section 4.4).

5 *Multivariate distributions*

5.5 Wishart distribution

The Wishart distribution is a multivariate generalisation of the gamma distribution.

Recall that the gamma distribution can be motivated as the distribution of sums of squared normal random variables. Likewise, the Wishart distribution can be understood as the sum of squared multivariate normal variables:

$$z_1, z_2, \ldots, z_k \stackrel{\mathrm{iid}}{\sim} N_d(0, S)$$

with $S = (s_{ij})$ the specified covariance matrix. The random variable

$$\underbrace{\mathbf{X}}_{d \times d} = \sum_{i=1}^{k} \underbrace{\mathbf{z}_i \mathbf{z}_i^T}_{d \times d}$$

is a random *matrix* and is distributed as

$$\boldsymbol{X} \sim \mathbf{W}_d\left(\boldsymbol{S}, \boldsymbol{k}\right)$$

with mean

$$E(X) = kS$$

and variances

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

We often also employ the Wishart distribution in **mean parametrisation**

$$W_d\left(S = \frac{M}{k}, k\right)$$

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

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

and variances are

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

If S or M is a scalar rather than a matrix then the multivariate Wishart distribution reduces to the univariate Wishart aka gamma distribution (Section 4.5).

5.6 Inverse Wishart distribution

The inverse Wishart distribution is a multivariate generalisation of the inverse gamma distribution and is linked to the Wishart distribution.

A random matrix **X** following an **inverse Wishart distribution** is denoted by

$$X \sim W_d^{-1}(\Psi, \nu)$$

where Ψ is the scale parameter and ν the shape parameter. The corresponding mean is given by

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

and the variances are

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

The inverse of *X* is then Wishart distributed:

$$X^{-1} \sim W_d \left(S = \Psi^{-1}, k = \nu \right)$$

Instead of Ψ and ν we may use the mean parametrisation with parameters $\kappa = \nu - d - 1$ and $M = \Psi/(\nu - d - 1)$:

$$\boldsymbol{X} \sim \mathrm{W}_d^{-1} \left(\boldsymbol{\Psi} = \kappa \boldsymbol{M} \right), \ \boldsymbol{\nu} = \kappa + d + 1$$

E(X) = M

with mean

and variances

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

If Ψ or M is a scalar rather than a matrix then the multivariate inverse Wishart distribution reduces to the univariate inverse Wishart aka inverse gamma distribution (Section 4.6).

Bibliography

Mardia, K. V., J. T. Kent, and J. M. Bibby. 1979. *Multivariate Analysis*. Academic Press.

Whittle, P. 2000. *Probability via Expectation*. 3rd ed. Springer. https://doi.org/10.1 007/978-1-4612-0509-8.