

Fundamentals of Probability Theory
- Notes -

Lim Zi Xiang

August 11, 2025

Contents

Chapter 1

Random variables and random vectors Page 2

- 1.1 Random vectors 2
Discrete random vectors — 3 • Continuous random vectors — 4 • Random vectors in general — 5 • Joint distribution — 5 • Random matrices — 5 • The marginal distribution of a random vector — 5 • Marginalisation of a joint distribution — 6 • The marginal distribution of a discrete vector — 6 • Marginalisation of a discrete distribution — 6 • The marginal distribution of a continuous vector — 6 • Marginalisation of a continuous distribution — 7 • Partial derivatives of the distribution function of a continuous vector — 7 • A more rigorous definition of random vectors — 7 • Exercises — 7
- 1.2 Expected value 8
Expected value of a discrete random variable — 8 • Expected value of a continuous random variable — 9 • Expected value of a random variable in general: the Riemann-Stieltjes integral — 9 • Expected value of a random variable in general: the Lebesgue integral (skipped) — 10 • The transformation theorem — 10 • Linearity of the expected value — 11 • Expected value of random vectors — 11 • Expected value of random matrices — 11 • Integrability and L_p spaces — 11 • Exercises — 11
- 1.3 Properties of the expected value 11
Scalar multiplication of a random variable — 11 • Sums of random variables — 11 • Linear combination of random variables — 11 • Expected value of a constant — 11 • Expectation of a product of random variables — 11 • Non-linear transformations — 11 • Addition of a constant matrix and a matrix with random entries — 11 • Multiplication of a constant matrix and a matrix with random entries — 11 • Expectation of a positive random variable — 11 • Preservation of almost sure inequalities — 11 • Exercises — 11
- 1.4 Variance 11
- 1.5 Covariance 11
- 1.6 Linear correlation 11
- 1.7 Covariance matrix 11
- 1.8 Indicator functions 11
- 1.9 Quantile 11

Chapter 2

Conditional distributions and independence Page 12

Rigorous conditional probability — 12

Chapter 1

Random variables and random vectors

1.1 Random vectors

Definition 1.1.1: Random vector

Let Ω be a sample space. A **random vector** \mathbf{X} is a function from the sample space Ω to the set of K -dimensional real vectors \mathbb{R}^K :

$$\mathbf{X} : \Omega \rightarrow \mathbb{R}^K.$$

To put it simply, a random vector is a vector whose value depends on the outcome of the probabilistic experiment. The real vector $\mathbf{X}(\omega)$ associated to a sample point $\omega \in \Omega$ is a **realisation** of the random vector. The set of all possible realisations is the **support**, denoted R_X .

Note:-

We denote the probability of an event $E \subseteq \Omega$ by $P(E)$. We use the following conventions when dealing with random vectors:

- For $A \subseteq \mathbb{R}^K$, $P(\mathbf{X} \in A) = P(\{\omega \in \Omega : \mathbf{X}(\omega) \in A\})$.
- For $A \subseteq \mathbb{R}^K$, $P_X(A) = P(\mathbf{X} \in A)$.

It is very common in applied work to build statistical models where a random vector \mathbf{X} is defined by directly specifying P_X and omitting the specification of the sample space Ω .

- We often write \mathbf{X} to mean $\mathbf{X}(\omega)$.

Example 1.1.1 (Defining a random vector on a sample space)

Two coins are tossed. The possible outcomes of each toss can be either tail (T) or head (H). The sample space is

$$\Omega = \{TT, TH, HT, HH\}.$$

The four possible outcomes are assigned equal probabilities:

$$P(\{TT\}) = P(\{TH\}) = P(\{HT\}) = P(\{HH\}) = \frac{1}{4}.$$

If the outcome is tails, we win a dollar, otherwise we lose one dollar. A 2D random vector \mathbf{X} indicates the amount we win on each toss:

$$\mathbf{X}(\omega) = \begin{cases} \begin{pmatrix} 1 & 1 \end{pmatrix} & \text{if } \omega = TT \\ \begin{pmatrix} 1 & -1 \end{pmatrix} & \text{if } \omega = TH \\ \begin{pmatrix} -1 & 1 \end{pmatrix} & \text{if } \omega = HT \\ \begin{pmatrix} -1 & -1 \end{pmatrix} & \text{if } \omega = HH \end{cases}$$

The probability of winning one dollar on both tosses is

$$P\left(\mathbf{X} = \begin{pmatrix} 1 & 1 \end{pmatrix}\right) = P\left(\{\omega \in \Omega : \mathbf{X}(\omega) = \begin{pmatrix} 1 & 1 \end{pmatrix}\}\right) = P(\{TT\}) = \frac{1}{4}.$$

The probability of losing one dollar on the second toss is

$$P(X_2 = -1) = P(\{\omega \in \Omega : X_2(\omega) = -1\}) = P(\{TH, HH\}) = \frac{1}{2}.$$

1.1.1 Discrete random vectors

Definition 1.1.2: Discrete random vector

A random vector \mathbf{X} is **discrete** iif

1. its support $R_{\mathbf{X}}$ is a countable set;
2. there is a function $p_{\mathbf{X}} : \mathbb{R}^K \rightarrow [0, 1]$, called the **joint probability mass function** of \mathbf{X} , such that for any $\mathbf{x} \in \mathbb{R}^K$:

$$p_{\mathbf{X}}(\mathbf{x}) = \begin{cases} P(\mathbf{X} = \mathbf{x}) & \text{if } \mathbf{x} \in R_{\mathbf{X}}; \\ 0 & \text{if } \mathbf{x} \notin R_{\mathbf{X}}. \end{cases}$$

Note:-

The following are equivalent notations used interchangeably to indicate the joint pmf:

$$p_{\mathbf{X}}(\mathbf{x}) = p_{\mathbf{X}}(x_1, \dots, x_K) = p_{X_1, \dots, X_K}(x_1, \dots, x_K).$$

Example 1.1.2

Suppose \mathbf{X} is a 2D random vector whose components (X_1 and X_2) can take only two values: 1 or 0, and the four possible combinations of 0 and 1 are equally likely. The support of the discrete vector \mathbf{X} is

$$R_{\mathbf{X}} = \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right\}.$$

Its pmf is

$$p_{\mathbf{X}} = \begin{cases} 0.25 & \text{if } \mathbf{x} = \begin{pmatrix} 1 & 1 \end{pmatrix}^T; \\ 0.25 & \text{if } \mathbf{x} = \begin{pmatrix} 1 & 0 \end{pmatrix}^T; \\ 0.25 & \text{if } \mathbf{x} = \begin{pmatrix} 0 & 1 \end{pmatrix}^T; \\ 0.25 & \text{if } \mathbf{x} = \begin{pmatrix} 0 & 0 \end{pmatrix}^T; \\ 0 & \text{otherwise.} \end{cases}$$

1.1.2 Continuous random vectors

Definition 1.1.3: Continuous random vector

A random vector \mathbf{X} is **continuous** (or **absolutely continuous**) iff

1. its support R_X is uncountable;
2. there exists a function $f_X : \mathbb{R}^K \rightarrow [0, \infty]$, called the **joint probability density function** of \mathbf{X} , such that for any set $A \subseteq \mathbb{R}^K$ where

$$A = [a_1, b_1] \times \dots \times [a_K, b_K].$$

The probability that $\mathbf{X} \in A$ is calculated by

$$P(\mathbf{X} \in A) = \int_{a_1}^{b_1} \dots \int_{a_K}^{b_K} f_X(x_1, \dots, x_K) dx_K \dots dx_1$$

provided the multiple integral is well defined.

Example 1.1.3

Suppose \mathbf{X} is a 2D random vector whose components X_1 and X_2 are independent uniform random variables on the interval $[0, 1]$. Then, \mathbf{X} is an example of a continuous vector with support

$$R_X = [0, 1] \times [0, 1].$$

Its joint pmf is

$$f_X(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in [0, 1] \times [0, 1]; \\ 0 & \text{otherwise.} \end{cases}$$

The probability that the realisation of \mathbf{X} falls in the rectangle $[0, 0.5] \times [0, 0.5]$ is

$$\begin{aligned} P(\mathbf{X} \in [0, 0.5] \times [0, 0.5]) &= \int_0^{0.5} \int_0^{0.5} f_X(x_1, x_2) dx_2 dx_1 \\ &= \int_0^{0.5} \int_0^{0.5} (1) dx_2 dx_1 \\ &= \int_0^{0.5} [x_2]_0^{0.5} dx_1 \\ &= \int_0^{0.5} 0.5 dx_1 \\ &= [0.5x_1]_0^{0.5} \\ &= 0.25 \end{aligned}$$

1.1.3 Random vectors in general

Definition 1.1.4: Joint distribution function

Let \mathbf{X} be a random vector. The **joint (cumulative) distribution function** of \mathbf{X} is a function $F_{\mathbf{X}} : \mathbb{R}^K \rightarrow [0, 1]$ such that

$$F_{\mathbf{X}}(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_K \leq x_K), \forall \mathbf{x} \in \mathbb{R}^K,$$

where the components of \mathbf{X} and \mathbf{x} are denoted by X_k and x_k respectively, for $k = 1, \dots, K$.

Similarly for the case of joint pmf/pdf, the following notations are used interchangeably to indicate the joint cdf:

$$F_{\mathbf{X}}(\mathbf{x}) = F_{\mathbf{X}}(x_1, \dots, x_K) = F_{X_1, \dots, X_K}(x_1, \dots, x_K).$$

1.1.4 Joint distribution

Sometimes we talk about the **joint distribution** of a random vector without specifying whether we mean the joint cdf, pmf, or pdf. And this is justified, since the joint pmf/pdf completely determines and is completely determined by the joint cdf of a discrete/continuous vector.

1.1.5 Random matrices

Definition 1.1.5: Random matrix

A random matrix is a matrix whose entries are random variables.

A random matrix can always be written as a random vector by vectorising it: given a $K \times L$ random matrix \mathbf{A} , its vectorisation, denoted $\text{vec}(\mathbf{A})$ is the $KL \times 1$ random vector obtained by stacking the columns of \mathbf{A} on top of each other.

Example 1.1.4

Let \mathbf{A} be the following 2×2 random matrix:

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}.$$

The vectorisation of \mathbf{A} is the following 4×1 vector:

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{pmatrix}.$$

If $\text{vec}(\mathbf{A})$ is a discrete/continuous vector, then \mathbf{A} is a **discrete/continuous random matrix**, the joint pmf of \mathbf{A} is just the joint pmf/pdf of $\text{vec}(\mathbf{A})$.

1.1.6 The marginal distribution of a random vector

Let X_i be the i -th component of a K -dimensional random vector \mathbf{X} . The cdf $F_{X_i}(\mathbf{x})$ of X_i is the marginal distribution function of X_i .

If \mathbf{X} is discrete/continuous, then X_i is a discrete/continuous random variable and its pmf $p_{X_i}(\mathbf{x})$ /pdf $f_{X_i}(\mathbf{x})$ is the **marginal pmf/pdf** of X_i .

1.1.7 Marginalisation of a joint distribution

Marginalisation is the process of deriving the distribution of a component X_i of a random vector \mathbf{X} from the joint distribution of \mathbf{X} .

It can also have a broader meaning of deriving the joint distribution of a subset of the set of components of \mathbf{X} from the joint distribution of \mathbf{X} , e.g. if \mathbf{X} has three components X_1, X_2, X_3 , we can marginalise their joint distribution to find the joint distribution of X_1 and X_2 . In this case, X_3 is said to be marginalised out of the joint distribution of X_1, X_2 , and X_3 .

1.1.8 The marginal distribution of a discrete vector

Let X_i be the i -th component of a K -dimensional discrete random vector \mathbf{X} . The marginal pmf of X_i is derived from the joint pmf by:

$$p_{X_i}(x) = \sum_{(x_1, \dots, x_K) \in R_{\mathbf{X}} : x_i = x} p_{\mathbf{X}}(x_1, \dots, x_K),$$

where the sum is over the set

$$\{(x_1, \dots, x_K) \in R_{\mathbf{X}} : x_i = x\},$$

i.e. the probability that $X_i = x$ is obtained as the sum of the probabilities of all the vectors in $R_{\mathbf{X}}$ such that their i -th component is equal to x .

1.1.9 Marginalisation of a discrete distribution

Let X_i be the i -th component of a K -dimensional discrete random vector \mathbf{X} . By marginalising X_i out of the joint distribution of \mathbf{X} , we obtain the joint distribution of the remaining components of \mathbf{X} , \mathbf{X}_{-i} :

$$\mathbf{X}_{-i} = (X_1 \quad \dots \quad X_{i-1} \quad X_{i+1} \quad \dots \quad X_K).$$

The joint pmf of \mathbf{X}_{-i} is computed as follows:

$$p_{\mathbf{X}_{-i}}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_K) = \sum_{x_i \in R_{X_i}} p_{\mathbf{X}}(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_K),$$

i.e. the joint pmf of \mathbf{X}_{-i} is computed by summing the joint pmf of \mathbf{X} over all values of x_i that belong to the support of X_i .

1.1.10 The marginal distribution of a continuous vector

Let X_i be the i -th component of a K -dimensional continuous random vector \mathbf{X} . The **marginal pdf** of X_i is derived from the joint pdf of \mathbf{X} by

$$f_{X_i}(x) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\mathbf{X}}(x_1, \dots, x_K) dx_K \dots dx_{i+1} dx_{i-1} dx_1,$$

i.e. the joint pdf evaluated at $x_i = x$ is integrated with respect to all variables except x_i .

1.1.11 Marginalisation of a continuous distribution

Let X_i be the i -th component of a continuous random vector \mathbf{X} . By marginalising X_i out of the joint distribution of \mathbf{X} , we obtain the joint distribution of the remaining components of \mathbf{X} , \mathbf{X}_{-i} .

The joint pdf of \mathbf{X}_{-i} is then computed by

$$f_{\mathbf{X}_{-i}}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_K) = \int_{-\infty}^{\infty} f_{\mathbf{X}}(x_1, \dots, x_K) dx_i,$$

i.e. the joint pdf of \mathbf{X}_{-i} is computed by integrating the joint pdf of \mathbf{X} with respect to x_i .

1.1.12 Partial derivatives of the distribution function of a continuous vector

We know that if \mathbf{X} is continuous, then

$$F_{\mathbf{X}}(\mathbf{x}) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_K} f_{\mathbf{X}}(t_1, \dots, t_K) dt_K \dots dt_1.$$

So by taking the K -th order cross-partial derivative with respect to x_1, \dots, x_K of both sides of the above equation, we get

$$\frac{\partial^K F_{\mathbf{X}}(\mathbf{x})}{\partial x_1 \dots \partial x_K} = f_{\mathbf{X}}(\mathbf{x}).$$

1.1.13 A more rigorous definition of random vectors

The following is a more rigorous definition of random vector using the formalism of measure theory. (I'll ignore this part for the time being.)

Definition 1.1.6

Let (Ω, \mathcal{F}, P) be a probability space. Let $\mathcal{B}(\mathbb{R}^K)$ be the Borel sigma-algebra of \mathbb{R}^K (i.e. the smallest sigma-algebra containing all open hyper-rectangles in \mathbb{R}^K). A function $\mathbf{X} : \Omega \rightarrow \mathbb{R}^K$ such that

$$\{\omega \in \Omega : \mathbf{X}(\omega) \in B\} \in \mathcal{F}$$

for any $B \in \mathcal{B}(\mathbb{R}^K)$ is said to be a random vector on Ω .

This definition ensures that the probability that the realisation of \mathbf{X} belongs to a set $B \in \mathcal{B}(\mathbb{R}^K)$ that can be defined as

$$P(\mathbf{X} \in B) := P(\{\omega \in \Omega : \mathbf{X}(\omega) \in B\})$$

because the set $\{\omega \in \Omega : \mathbf{X}(\omega) \in B\}$ belongs to the sigma-algebra \mathcal{F} and hence its probability is well-defined.

1.1.14 Exercises

Question 1

Let \mathbf{X} be a 2×1 discrete random vector and denote its components by X_1 and X_2 .

Let the support of \mathbf{X} be the set of all 2×1 vectors such that their entries belong to the set of the first three natural numbers, that is,

$$R_{\mathbf{X}} = \{\mathbf{x} = \begin{pmatrix} x_1 & x_2 \end{pmatrix}^T : x_1 \in N_3 \text{ and } x_2 \in N_3\}.$$

1.2 Expected value

In this section we give an informal definition of expected value. A formal definition involves the Lebesgue integral which I will ignore for the time being.

Definition 1.2.1: Expected value

The **expected value** of a random variable \mathbf{X} , $E[X]$ is the weighted average of the values that \mathbf{X} can take on, where each possible value is weighted by its respective probability.

1.2.1 Expected value of a discrete random variable

Definition 1.2.2: Expected value of a discrete random variable

Let X be a discrete random variable with support R_X and pmf $p_X(x)$. The expected value of \mathbf{X} is

$$E[X] = \sum_{x \in R_X} x p_X(x),$$

provided that we have **absolute summability**

$$\sum_{x \in R_X} |x| p_X(x) < \infty,$$

ensuring that the summation is well-defined when R_X contained infinitely many elements.

Note:-

When summing infinitely many terms, the order in which you sum them can change the results and the expected value of \mathbf{X} is not well-defined or does not exist. However this is not true if the terms are absolutely summable.

Example 1.2.1 (Expected value)

Let \mathbf{X} be a random variable with support $R_X = \{0, 1\}$ and pmf

$$p_X(x) = \begin{cases} 0.5 & \text{if } x = 1; \\ 0.5 & \text{if } x = 0; \\ 0 & \text{otherwise.} \end{cases}$$

Its expected value is

$$\begin{aligned} E[X] &= \sum_{x \in R_X} x p_X(x) \\ &= (1)(0.5) + (0)(0.5) \\ &= 0.5 \end{aligned}$$

1.2.2 Expected value of a continuous random variable

Definition 1.2.3: Expected value of a continuous random variable

Let X be a continuous random variable with pdf $f_X(x)$. The expected value of X is

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx,$$

provided that we have absolute integrability

$$\int_{-\infty}^{\infty} |x| f_X(x) dx < \infty.$$

Example 1.2.2 (Expected value of continuous random variable)

Let X be a continuous random variable with support $R_X = [0, \infty)$ and pdf

$$f_X(x) = \begin{cases} \lambda \exp(-\lambda x) & \text{if } x \in [0, \infty); \\ 0 & \text{otherwise.} \end{cases}$$

where $\lambda > 0$. Its expected value is

$$\begin{aligned} E[X] &= \int_{-\infty}^{\infty} x f_X(x) dx \\ &= \int_0^{\infty} x \lambda (-\lambda x) dx \\ &= \frac{1}{\lambda} \int_{t=0}^{t=\infty} t \exp(-t) dt \\ &= \frac{1}{\lambda} \left\{ [-t \exp(-t)]_{t=0}^{t=\infty} + \int_0^{\infty} \exp(-t) dt \right\} \end{aligned}$$

1.2.3 Expected value of a random variable in general: the Riemann-Stieltjes integral

Definition 1.2.4: Expected value (general)

Let X be a random variable with cdf $F_X(x)$. The expected value of X is

$$E[X] = \int_{-\infty}^{\infty} x dF_X(x),$$

where the integral is a Riemann-Stieltjes integral and the expected value exists and is well-defined iff the integral is well-defined.

This definition gives a formal notation which allows for a unified treatment of discrete and continuous random variables and can be treated as a sum in one case and as ordinary Riemann integral in the other.

Example 1.2.3

Let X be a random variable with support $R_X = [0, 1]$ and distribution function

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 0.5x & \text{if } 0 \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

Its expected value is

$$\begin{aligned} E[X] &= \int_{-\infty}^{\infty} x dF_X(x) \\ &= \int_0^1 x dF_X(x) + 1 \cdot \left[F_X(1) - \lim_{x \rightarrow 1^+} F_X(x) \right] \\ &= \int_0^1 x \frac{d}{dx} \left(\frac{1}{2}x \right) dx + 1 \cdot \left[1 - \frac{1}{2} \right] \\ &= \left[\frac{1}{4}x^2 \right]_0^1 + \frac{1}{2} \\ &= \frac{3}{4} \end{aligned}$$

1.2.4 Expected value of a random variable in general: the Lebesgue integral

Definition 1.2.5: Expected value (rigorous)

Let Ω be a sample space, P a probability measure defined on the events of Ω and X a random variable defined on Ω . The expected value of X is

$$E[X] = \int X dP,$$

provided the Lebesgue integral of X with respect to P exists and is well-defined.

1.2.5 The transformation theorem

Let X be a random variable, $g : \mathbb{R} \rightarrow \mathbb{R}$ be a real function. Define a new random variable Y as $Y = g(X)$. Then,

$$E[Y] = \int_{-\infty}^{\infty} g(x) dF_X(x)$$

provided that the integral exists. For discrete random variables, the formula becomes

$$E[Y] = \sum_{x \in R_X} g(x) p_X(x)$$

while for continuous random variables,

$$E[Y] = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

When \mathbf{X} is a discrete random vector and $p_{\mathbf{X}}(x)$ is its joint pmf, then

$$E[Y] = \sum_{x \in R_{\mathbf{X}}} g(x) p_{\mathbf{X}}(x).$$

When \mathbf{X} is a continuous random vector and $f_{\mathbf{X}}(x)$ is its joint pdf, then

$$E[Y] = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x) f_X(x) dx_1 \dots dx_K.$$

1.2.6 Linearity of the expected value

1.2.7 Expected value of random vectors

1.2.8 Expected value of random matrices

1.2.9 Integrability and Lp spaces

1.2.10 Exercises

1.3 Properties of the expected value

1.3.1 Scalar multiplication of a random variable

1.3.2 Sums of random variables

1.3.3 Linear combination of random variables

1.3.4 Expected value of a constant

1.3.5 Expectation of a product of random variables

1.3.6 Non-linear transformations

1.3.7 Addition of a constant matrix and a matrix with random entries

1.3.8 Multiplication of a constant matrix and a matrix with random entries

1.3.9 Expectation of a positive random variable

1.3.10 Preservation of almost sure inequalities

1.3.11 Exercises

1.4 Variance

1.5 Covariance

1.6 Linear correlation

1.7 Covariance matrix

1.8 Indicator functions

1.9 Quantile

Chapter 2

Conditional distributions and independence

2.0.1 Rigorous conditional probability