

Appetizer

Using Probability to Cover a Geometric Set

We begin our study of high-dimensional probability with an elegant argument that showcases the usefulness of probabilistic reasoning in geometry.

Recall that a *convex combination* of points $z_1, \dots, z_m \in \mathbb{R}^n$ is a linear combination with coefficients that are non-negative and sum to 1, i.e., it is a sum of the form

$$\sum_{i=1}^m \lambda_i z_i \quad \text{where } \lambda_i \geq 0 \quad \text{and} \quad \sum_{i=1}^m \lambda_i = 1. \quad (0.1)$$

The *convex hull* of a set $T \subset \mathbb{R}^n$ is the set of all convex combinations of all finite collections of points in T :

$$\text{conv}(T) := \{\text{convex combinations of } z_1, \dots, z_m \in T \text{ for } m \in \mathbb{N}\};$$

see Figure 0.1 for illustration.

The number m of elements defining a convex combination in \mathbb{R}^n is not restricted a priori. However, the classical theorem of Caratheodory states that one can always take $m \leq n + 1$.

Theorem 0.0.1 (Caratheodory's theorem) *Every point in the convex hull of a set $T \subset \mathbb{R}^n$ can be expressed as a convex combination of at most $n + 1$ points from T .*

The bound $n + 1$ cannot be improved, as it is clearly attained for a simplex T (a set of $n + 1$ points in general positions). Suppose, however, that we want only to *approximate* a point $x \in \text{conv}(T)$ rather than to represent it exactly as a convex combination. Can we do

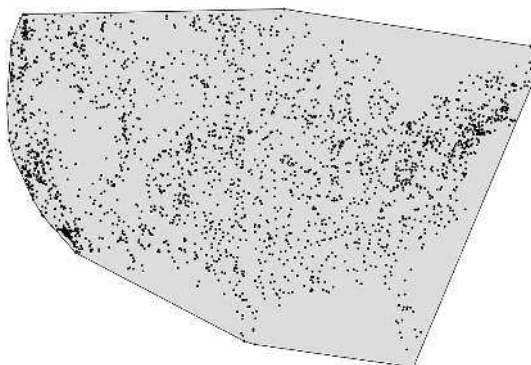


Figure 0.1 The convex hull of a set of points representing major US cities.

this with fewer than $n + 1$ points? We now show that it is possible, and actually the number of required points does not need to depend on the dimension n at all!

Theorem 0.0.2 (Approximate form of Caratheodory’s theorem) *Consider a set $T \subset \mathbb{R}^n$ whose diameter¹ is bounded by 1. Then, for every point $x \in \text{conv}(T)$ and every integer k , one can find points $x_1, \dots, x_k \in T$ such that*

$$\left\| x - \frac{1}{k} \sum_{j=1}^k x_j \right\|_2 \leq \frac{1}{\sqrt{k}}.$$

There are two reasons why this result is surprising. First, the number of points k in convex combinations does not depend on the dimension n . Second, the coefficients of convex combinations can be made all equal. (Note, however, that repetitions among the points x_i are allowed.)

Proof Our argument is known as the *empirical method* of B. Maurey.

Translating T if necessary, we may assume that not only the diameter but also the *radius* of T is bounded by 1, i.e.,

$$\|t\|_2 \leq 1 \quad \text{for all } t \in T. \tag{0.2}$$

Fix a point $x \in \text{conv}(T)$ and express it as a convex combination of some vectors $z_1, \dots, z_m \in T$ as in (0.1). Now, interpret the definition of the convex combination (0.1) probabilistically, with the λ_i taking the roles of probabilities. Specifically, we can define a random vector Z that takes the value z_i with probability λ_i :

$$\mathbb{P}\{Z = z_i\} = \lambda_i, \quad i = 1, \dots, m.$$

(This is possible by the fact that the weights λ_i are non-negative and sum to 1.) Then

$$\mathbb{E} Z = \sum_{i=1}^m \lambda_i z_i = x.$$

Consider independent copies Z_1, Z_2, \dots of Z . By the strong law of large numbers,

$$\frac{1}{k} \sum_{j=1}^k Z_j \rightarrow x \quad \text{almost surely as } k \rightarrow \infty.$$

To get a quantitative form of this result, let us compute the variance of $\frac{1}{k} \sum_{j=1}^k Z_j$. (Incidentally, this computation is at the heart of the proof of the weak law of large numbers.) We obtain

$$\begin{aligned} \mathbb{E} \left\| x - \frac{1}{k} \sum_{j=1}^k Z_j \right\|_2^2 &= \frac{1}{k^2} \mathbb{E} \left\| \sum_{j=1}^k (Z_j - x) \right\|_2^2 \quad (\text{since } \mathbb{E}(Z_i - x) = 0) \\ &= \frac{1}{k^2} \sum_{j=1}^k \mathbb{E} \|Z_j - x\|_2^2. \end{aligned}$$

¹ The diameter of T is defined as $\text{diam}(T) = \sup\{\|s - t\|_2 : s, t \in T\}$. We have assumed that $\text{diam}(T) = 1$ for simplicity. For a general set T , the bound in the theorem changes to $\text{diam}(T)/\sqrt{k}$. Check this!

The last identity is just a higher-dimensional version of the basic fact that the variance of a sum of independent random variables equals the sum of the variances; see Exercise 0.0.3 below.

It remains to bound the variances of the terms. We have

$$\begin{aligned} \mathbb{E} \|Z_j - x\|_2^2 &= \mathbb{E} \|Z - \mathbb{E} Z\|_2^2 \\ &= \mathbb{E} \|Z\|_2^2 - \|\mathbb{E} Z\|_2^2 \quad (\text{another variance identity; see Exercise 0.0.3}) \\ &\leq \mathbb{E} \|Z\|_2^2 \leq 1 \quad (\text{since } Z \in T \text{ and using (0.2)}). \end{aligned}$$

We have shown that

$$\mathbb{E} \left\| x - \frac{1}{k} \sum_{j=1}^k Z_j \right\|_2^2 \leq \frac{1}{k}.$$

Therefore, there exists a realization of the random variables Z_1, \dots, Z_k such that

$$\left\| x - \frac{1}{k} \sum_{j=1}^k Z_j \right\|_2^2 \leq \frac{1}{k}.$$

Since by construction each Z_j takes values in T , the proof is complete. ■

Exercise 0.0.3 ☞☞ Check the following variance identities, which we used in the proof of Theorem 0.0.2.

(a) Let Z_1, \dots, Z_k be independent mean-zero random vectors in \mathbb{R}^n . Show that

$$\mathbb{E} \left\| \sum_{j=1}^k Z_j \right\|_2^2 = \sum_{j=1}^k \mathbb{E} \|Z_j\|_2^2.$$

(b) Let Z be a random vector in \mathbb{R}^n . Show that

$$\mathbb{E} \|Z - \mathbb{E} Z\|_2^2 = \mathbb{E} \|Z\|_2^2 - \|\mathbb{E} Z\|_2^2.$$

Let us give one application of Theorem 0.0.2 in computational geometry. Suppose that we are given a subset $P \subset \mathbb{R}^n$ and asked to cover it by balls of a given radius ε ; see Figure 0.2. What is the smallest number of balls needed, and how should we place them?

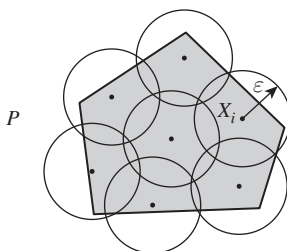


Figure 0.2 The covering problem asks how many balls of radius ε are needed to cover a given set P in \mathbb{R}^n and where to place these balls.

Corollary 0.0.4 (Covering polytopes by balls) *Let P be a polytope in \mathbb{R}^n with N vertices and whose diameter is bounded by 1. Then P can be covered by at most $N^{\lceil 1/\varepsilon^2 \rceil}$ Euclidean balls of radii $\varepsilon > 0$.*

Proof Let us define the centers of the balls as follows. Let $k := \lceil 1/\varepsilon^2 \rceil$ and consider the set

$$\mathcal{N} := \left\{ \frac{1}{k} \sum_{j=1}^k x_j : x_j \text{ are vertices of } P \right\}.$$

We claim that the family of ε -balls centered at \mathcal{N} satisfies the conclusion of the corollary. To check this, note that the polytope P is the convex hull of the set of its vertices, which we denote by T . Thus we can apply Theorem 0.0.2 to any point $x \in P = \text{conv}(T)$ and deduce that x is within a distance $1/\sqrt{k} \leq \varepsilon$ from some point in \mathcal{N} . This shows that the ε -balls centered at \mathcal{N} do indeed cover P .

To bound the cardinality of \mathcal{N} , note that there are N^k ways to choose k out of N vertices with repetition. Thus $|\mathcal{N}| \leq N^k = N^{\lceil 1/\varepsilon^2 \rceil}$. The proof is complete. ■

In this book we will learn several other approaches to the covering problem in relation to packing (Section 4.2), entropy and coding (Section 4.3), and random processes (Chapters 7 and 8).

To finish this section, let us show how to slightly improve Corollary 0.0.4.

Exercise 0.0.5 (The sum of binomial coefficients) 🍷🍷 Prove the inequalities

$$\left(\frac{n}{m}\right)^m \leq \binom{n}{m} \leq \sum_{k=0}^m \binom{n}{k} \leq \left(\frac{en}{m}\right)^m$$

for all integers $m \in [1, n]$. ■

Exercise 0.0.6 (Improved covering) 🍷🍷 Check that, in Corollary 0.0.4,

$$(C + C\varepsilon^2 N)^{\lceil 1/\varepsilon^2 \rceil}$$

Euclidean balls suffice. Here C is a suitable absolute constant. (Note that this bound is slightly stronger than $N^{\lceil 1/\varepsilon^2 \rceil}$ for small ε .) ■

0.1 Notes

In this appetizer we gave an illustration of the *probabilistic method*, where one employs randomness to construct a useful object. The book [8] presents many illustrations of the probabilistic method, mainly in combinatorics.

The empirical method of B. Maurey presented in this section was originally proposed in [162]. B. Carl used it to get bounds on covering numbers [48] including those stated in Corollary 0.0.4 and Exercise 0.0.6. The bound in Exercise 0.0.6 is sharp [48, 49].

1

Preliminaries on Random Variables

In this chapter we recall some basic concepts and results of probability theory. The reader should already be familiar with most of this material, which is routinely taught in introductory probability courses.

Expectation, variance, and moments of random variables are introduced in Section 1.1. Some classical inequalities can be found in Section 1.2. The two fundamental limit theorems of probability – the law of large numbers and the central limit theorem – are recalled in Section 1.3.

1.1 Basic Quantities Associated with Random Variables

In basic courses in probability theory, one learns about the two most important quantities associated with a random variable X , namely the *expectation*¹ (also called the *mean*) and *variance*. They will be denoted in this book by²

$$\mathbb{E} X \quad \text{and} \quad \text{Var}(X) = \mathbb{E}(X - \mathbb{E} X)^2.$$

Let us recall some other classical quantities and functions that describe probability distributions. The *moment generating function* of X is defined as

$$M_X(t) = \mathbb{E} e^{tX}, \quad t \in \mathbb{R}.$$

For $p > 0$, the *p th moment* of X is defined as $\mathbb{E} X^p$, and the *p th absolute moment* is $\mathbb{E} |X|^p$.

It is useful to take the p th root of the moments, which leads to the notion of the L^p norm of a random variable:

$$\|X\|_{L^p} = (\mathbb{E} |X|^p)^{1/p}, \quad p \in (0, \infty).$$

This definition can be extended to $p = \infty$ by the essential supremum of $|X|$:

$$\|X\|_{L^\infty} = \text{ess sup } |X|.$$

For fixed p and a given probability space $(\Omega, \Sigma, \mathbb{P})$, the classical vector space $L^p = L^p(\Omega, \Sigma, \mathbb{P})$ consists of all random variables X on Ω with finite L^p norm, that is,

¹ If you have studied measure theory, you will recall that the expectation $\mathbb{E} X$ of a random variable X on a probability space $(\Omega, \Sigma, \mathbb{P})$ is, by definition, the Lebesgue integral of the function $X: \Omega \rightarrow \mathbb{R}$. This makes all theorems on Lebesgue integration applicable in probability theory for expectations of random variables.

² Throughout this book, we omit brackets and simply write $\mathbb{E} f(X)$. Thus, nonlinear functions bind before an expectation.

$$L^p = \{X: \|X\|_{L^p} < \infty\}.$$

If $p \in [1, \infty]$, the quantity $\|X\|_{L^p}$ is a norm and L^p is a *Banach space*. This fact follows from Minkowski's inequality, which we recall in (1.4). For $p < 1$, the triangle inequality fails and $\|X\|_{L^p}$ is not a norm.

The exponent $p = 2$ is special in that L^2 is not only a Banach space but also a *Hilbert space*. The inner product and the corresponding norm on L^2 are given by

$$\langle X, Y \rangle_{L^2} = \mathbb{E} XY, \quad \|X\|_{L^2} = (\mathbb{E} |X|^2)^{1/2}. \quad (1.1)$$

Then the *standard deviation* of X can be expressed as

$$\|X - \mathbb{E} X\|_{L^2} = \sqrt{\text{Var}(X)} = \sigma(X).$$

Similarly, we can express the *covariance* of random variables X and Y as

$$\text{cov}(X, Y) = \mathbb{E}((X - \mathbb{E} X)(Y - \mathbb{E} Y)) = \langle X - \mathbb{E} X, Y - \mathbb{E} Y \rangle_{L^2}. \quad (1.2)$$

Remark 1.1.1 (Geometry of random variables) When we consider random variables as vectors in the Hilbert space L^2 , the identity (1.2) gives a *geometric interpretation* of the notion of covariance: the more the vectors $X - \mathbb{E} X$ and $Y - \mathbb{E} Y$ are aligned with each other, the larger are their inner product and covariance.

1.2 Some Classical Inequalities

Jensen's inequality states that for any random variable X and a *convex*³ function $\varphi: \mathbb{R} \rightarrow \mathbb{R}$, we have

$$\varphi(\mathbb{E} X) \leq \mathbb{E} \varphi(X).$$

As a simple consequence of Jensen's inequality, $\|X\|_{L^p}$ is an *increasing function in p* , that is

$$\|X\|_{L^p} \leq \|X\|_{L^q} \quad \text{for any } 0 \leq p \leq q = \infty. \quad (1.3)$$

This inequality follows since $\phi(x) = x^{q/p}$ is a convex function if $q/p \geq 1$.

Minkowski's inequality states that for any $p \in [1, \infty]$ and any random variables $X, Y \in L^p$, we have

$$\|X + Y\|_{L^p} \leq \|X\|_{L^p} + \|Y\|_{L^p}. \quad (1.4)$$

This can be viewed as the *triangle inequality*, which implies that $\|\cdot\|_{L^p}$ is a norm when $p \in [1, \infty]$.

The *Cauchy–Schwarz inequality* states that, for any random variables $X, Y \in L^2$, we have

$$|\mathbb{E} XY| \leq \|X\|_{L^2} \|Y\|_{L^2}.$$

The more general *Hölder's inequality* states that if $p, q \in (1, \infty)$ are conjugate exponents, that is, $1/p + 1/q = 1$, then the random variables $X \in L^p$ and $Y \in L^q$ satisfy

³ By definition, a function φ is *convex* if $\varphi(\lambda x + (1 - \lambda)y) \leq \lambda\varphi(x) + (1 - \lambda)\varphi(y)$ for all $\lambda \in [0, 1]$ and all vectors x, y in the domain of φ .

$$|\mathbb{E} XY| \leq \|X\|_{L^p} \|Y\|_{L^q}.$$

This inequality also holds for the pair $p = 1, q = \infty$.

As we recall from basic probability concepts, the *distribution* of a random variable X is, intuitively, the information about what values X takes with what probabilities. More rigorously, the distribution of X is determined by the *cumulative distribution function* (CDF) of X , defined as

$$F_X(t) = \mathbb{P}\{X \leq t\}, \quad t \in \mathbb{R}.$$

It is often more convenient to work with the *tails* of random variables, namely with

$$\mathbb{P}\{X > t\} = 1 - F_X(t).$$

There is an important connection between the tails and the expectation (and more generally, the moments) of a random variable. The following identity is typically used to bound the expectation by the tails.

Lemma 1.2.1 (Integral identity) *Let X be a non-negative random variable X . Then*

$$\mathbb{E} X = \int_0^\infty \mathbb{P}\{X > t\} dt.$$

The two sides of this identity are either finite or infinite simultaneously.

Proof We can represent any non-negative real number x via the identity⁴

$$x = \int_0^x 1 dt = \int_0^\infty \mathbf{1}_{\{t < x\}} dt.$$

Substitute the random variable X for x and take expectation of both sides. This gives

$$\mathbb{E} X = \mathbb{E} \int_0^\infty \mathbf{1}_{\{t < X\}} dt = \int_0^\infty \mathbb{E} \mathbf{1}_{\{t < X\}} dt = \int_0^\infty \mathbb{P}\{t < X\} dt.$$

To change the order of expectation and integration in the second equality, we used the Fubini–Tonelli theorem. The proof is complete. ■

Exercise 1.2.2 (Generalization of integral identity)[♣] Prove the following extension of Lemma 1.2.1, which is valid for any random variable X (not necessarily non-negative):

$$\mathbb{E} X = \int_0^\infty \mathbb{P}\{X > t\} dt - \int_{-\infty}^0 \mathbb{P}\{X < t\} dt.$$

Exercise 1.2.3 (p th moment via the tail)[♣] Let X be a random variable and $p \in (0, \infty)$. Show that

$$\mathbb{E} |X|^p = \int_0^\infty pt^{p-1} \mathbb{P}\{|X| > t\} dt$$

whenever the right-hand side is finite. ■[♣]

⁴ Here and later in this book, $\mathbf{1}_E$ denotes the *indicator* of the event E ; it is the function that takes the value 1 if E occurs and 0 otherwise.

Another classical tool, Markov's inequality, can be used to bound the tail in terms of the expectation.

Proposition 1.2.4 (Markov's inequality) *For any non-negative random variable X and $t > 0$, we have*

$$\mathbb{P}\{X \geq t\} \leq \frac{\mathbb{E} X}{t}.$$

Proof Fix $t > 0$. We can represent any real number x via the identity

$$x = x\mathbf{1}_{\{x \geq t\}} + x\mathbf{1}_{\{x < t\}}.$$

Substitute the random variable X for x and take the expectation of both sides. This gives


$$\begin{aligned} \mathbb{E} X &= \mathbb{E} X\mathbf{1}_{\{X \geq t\}} + \mathbb{E} X\mathbf{1}_{\{X < t\}} \\ &\geq \mathbb{E} t\mathbf{1}_{\{X \geq t\}} + 0 = t \mathbb{P}\{X \geq t\}. \end{aligned}$$

Dividing both sides by t , we complete the proof. ■

A well-known consequence of Markov's inequality is Chebyshev's inequality. It offers a better, quadratic, dependence on t and, instead of controlling a one-side tail, it quantifies the *concentration* of X about its mean.

Corollary 1.2.5 (Chebyshev's inequality) *Let X be a random variable with mean μ and variance σ^2 . Then, for any $t > 0$, we have*

$$\mathbb{P}\{|X - \mu| \geq t\} \leq \frac{\sigma^2}{t^2}.$$

Exercise 1.2.6  Deduce Chebyshev's inequality by squaring both sides of the bound $|X - \mu| \geq t$ and applying Markov's inequality.

Remark 1.2.7 In Proposition 2.5.2 we will establish relations among the three basic quantities associated with random variables – the moment generating functions, the L^p norms, and the tails.

1.3 Limit Theorems

The study of *sums of independent random variables* is a core part of classical probability theory. Recall that the identity

$$\text{Var}(X_1 + \cdots + X_N) = \text{Var}(X_1) + \cdots + \text{Var}(X_N)$$

holds for any independent random variables X_1, \dots, X_N . If, furthermore, the X_i each have the same distribution, with mean μ and variance σ^2 , then dividing both sides by N we see that

$$\text{Var}\left(\frac{1}{N} \sum_{i=1}^N X_i\right) = \frac{\sigma^2}{N}. \quad (1.5)$$

Thus, the variance of the *sample mean* $\frac{1}{N} \sum_{i=1}^N X_i$ of the sample $\{X_1, \dots, X_N\}$ shrinks to zero as $N \rightarrow \infty$. This indicates that, for large N , we should expect that the sample mean concentrates tightly about its expectation μ . One of the most important results in probability theory – the law of large numbers – states precisely this.

Theorem 1.3.1 (Strong law of large numbers) *Let X_1, X_2, \dots be a sequence of i.i.d. random variables with mean μ . Consider the sum*

$$S_N = X_1 + \dots + X_N.$$

Then, as $N \rightarrow \infty$,

$$\frac{S_N}{N} \rightarrow \mu \quad \text{almost surely.}$$

The next result, the central limit theorem, goes one step further. It identifies the limiting distribution of the (properly scaled) sum of the X_i as the *normal* distribution, also called the *Gaussian* distribution. Recall that the *standard normal* distribution, denoted $N(0, 1)$, has density

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad x \in \mathbb{R}. \quad (1.6)$$

Theorem 1.3.2 (Lindeberg–Lévy central limit theorem) *Let X_1, X_2, \dots be a sequence of i.i.d. random variables with mean μ and variance σ^2 . Consider the sum*

$$S_N = X_1 + \dots + X_N$$

and normalize it to obtain a random variable with zero mean and unit variance as follows:

$$Z_N := \frac{S_N - \mathbb{E} S_N}{\sqrt{\text{Var}(S_N)}} = \frac{1}{\sigma\sqrt{N}} \sum_{i=1}^N (X_i - \mu).$$


Then, as $N \rightarrow \infty$,

$$Z_N \rightarrow N(0, 1) \quad \text{in distribution.}$$

Convergence in distribution means that the CDF of the normalized sum converges pointwise to the CDF of the standard normal distribution. We can express this in terms of tails as follows. Thus, for every $t \in \mathbb{R}$ we have

$$\mathbb{P}\{Z_N \geq t\} \rightarrow \mathbb{P}\{g \geq t\} = \frac{1}{\sqrt{2\pi}} \int_t^\infty e^{-x^2/2} dx$$

as $N \rightarrow \infty$, where $g \sim N(0, 1)$ is a standard normal random variable.

Exercise 1.3.3  Let X_1, X_2, \dots be a sequence of i.i.d. random variables with mean μ and finite variance. Show that

$$\mathbb{E} \left| \frac{1}{N} \sum_{i=1}^N X_i - \mu \right| = O\left(\frac{1}{\sqrt{N}}\right) \quad \text{as } N \rightarrow \infty.$$

One remarkable special case of the central limit theorem occurs when the X_i are Bernoulli random variables with some fixed parameter $p \in (0, 1)$, denoted

$$X_i \sim \text{Ber}(p).$$

Recall that this means that the X_i take the values 1 and 0 with probabilities p and $1 - p$ respectively; also recall that $\mathbb{E} X_i = p$ and $\text{Var}(X_i) = p(1 - p)$. The sum

$$S_N := X_1 + \cdots + X_N$$

is said to have the *binomial distribution* $\text{Binom}(N, p)$. The central limit theorem (Theorem 1.3.2) yields that, as $N \rightarrow \infty$,

$$\frac{S_N - Np}{\sqrt{Np(1-p)}} \rightarrow N(0, 1) \quad \text{in distribution.} \quad (1.7)$$

This special case of the central limit theorem is called the *de Moivre–Laplace theorem*.

Now suppose that $X_i \sim \text{Ber}(p_i)$, with parameters p_i that *decay to zero* as $N \rightarrow \infty$ so fast that the sum S_N has mean $O(1)$ instead of being proportional to N . The central limit theorem fails in this regime. A different result, which we are about to state, says that S_N still converges but to the *Poisson* instead of the normal distribution.

Recall that a random variable Z has a *Poisson distribution* with parameter λ , denoted

$$Z \sim \text{Pois}(\lambda),$$

if it takes values in $\{0, 1, 2, \dots\}$ with probabilities

$$\mathbb{P}\{Z = k\} = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k = 0, 1, 2, \dots \quad (1.8)$$

Theorem 1.3.4 (Poisson limit theorem) *Let $X_{N,i}$, $1 \leq i \leq N$, be independent random variables $X_{N,i} \sim \text{Ber}(p_{N,i})$, and let $S_N = \sum_{i=1}^N X_{N,i}$. Assume that, as $N \rightarrow \infty$,*

$$\max_{i \leq N} p_{N,i} \rightarrow 0 \quad \text{and} \quad \mathbb{E} S_N = \sum_{i=1}^N p_{N,i} \rightarrow \lambda < \infty.$$

Then, as $N \rightarrow \infty$,

$$S_N \rightarrow \text{Pois}(\lambda) \quad \text{in distribution.}$$

1.4 Notes

The material presented in this chapter is included in most graduate probability textbooks. In particular, proofs of the strong law of large numbers (Theorem 1.3.1) and the Lindeberg–Lévy central limit theorem (Theorem 1.3.2) can be found e.g. in [70, Sections 1.7 and 2.4] and [22, Sections 6 and 27].