

# The weak and strong laws of large numbers

Jordan Bell

jordan.bell@gmail.com

Department of Mathematics, University of Toronto

May 29, 2015

## 1 Introduction

Using Egorov's theorem, one proves that if a sequence of random variables  $X_n$  **converges almost surely** to a random variable  $X$  then  $X_n$  **converges in probability** to  $X$ .<sup>1</sup>

Let  $X_n$  be a sequence of  $L^1$  random variables,  $n \geq 1$ . **A weak law of large numbers** is a statement that

$$\frac{1}{n} \sum_{k=1}^n (X_k - E(X_k)) \tag{1}$$

converges in probability to 0. **A strong law of large numbers** is a statement that (1) converges almost surely to 0. Thus, if the hypotheses assumed on the sequence of random variables are the same, a strong law implies a weak law.

We shall prove the weak law of large numbers for a sequence of independent identically distributed  $L^1$  random variables, and the strong law of large for the same hypotheses. We give separate proofs for these theorems as an occasion to inspect different machinery, although to establish the weak law it thus suffices to prove the strong law. One reason to distinguish these laws is for cases when we impose different hypotheses.<sup>2</sup>

We also prove **Markov's weak law of large numbers**, which states that if  $X_n$  is a sequence of  $L^2$  random variables that are pairwise uncorrelated and

$$\frac{1}{n^2} \sum_{n=k}^n \text{Var}(X_k) \rightarrow 0,$$

then  $\frac{1}{n} \sum_{k=1}^n (X_k - E(X_k))$  converges to 0 in  $L^2$ , from which it follows using Chebyshev's inequality that it converges in probability to 0. (We remark that

---

<sup>1</sup> V. I. Bogachev, *Measure Theory*, volume I, p. 111, Theorem 2.2.3; <http://individual.utoronto.ca/jordanbell/notes/L0.pdf>, p. 3, Theorem 3.

<sup>2</sup>cf. Jordan M. Stoyanov, *Counterexamples in Probability*, third ed., p. 163, §15.3; Dieter Landers and Lothar Rogge, *Identically distributed uncorrelated random variables not fulfilling the WLLN*, Bull. Korean Math. Soc. **38** (2001), no. 3, 605–610.

a sequence of  $L^2$  random variables converging in  $L^2$  to 0 does not imply that it converges almost surely to 0, although there is indeed a subsequence that converges almost surely to 0.<sup>3</sup>)

If  $(\Omega, \mathcal{F}, P)$  is a probability space,  $(Y, \mathcal{A})$  is a measurable space, and  $T : (\Omega, \mathcal{F}) \rightarrow (Y, \mathcal{A})$  is measurable, the **pushforward measure of  $P$  by  $T$**  is

$$(T_*P)(A) = P(T^{-1}(A)), \quad A \in \mathcal{A}.$$

Then  $(Y, \mathcal{A}, T_*P)$  is a probability space. We remind ourselves of the **change of variables theorem**, which we shall use.<sup>4</sup> Let  $f : Y \rightarrow \mathbb{R}$  be a function. On the one hand, if  $f \in L^1(T_*P)$  then  $f \circ T \in L^1(P)$  and

$$\int_Y f d(T_*P) = \int_{\Omega} f \circ T dP. \quad (2)$$

On the other hand, if  $f$  is  $T_*P$ -measurable and  $f \circ T \in L^1(P)$ , then  $f \in L^1(T_*P)$  and

$$\int_Y f d(T_*P) = \int_{\Omega} f \circ T dP.$$

## 2 The weak law of large numbers

Suppose that  $X_n : (\Omega, \mathcal{F}, P) \rightarrow \mathbb{R}$ ,  $n \geq 1$ , are independent identically distributed  $L^1$  random variables, and write

$$S_n = \sum_{k=1}^n X_k,$$

for which  $E(S_n) = \sum_{k=1}^n E(X_k) = nE(X_1)$ .

**Lemma 1.** *If  $X_n$  are  $L^2$ , then for any  $\lambda > 0$  and for any  $n \geq 1$ ,*

$$P\left(\left|\frac{S_n}{n} - E(X_1)\right| \geq \lambda\right) \leq \frac{E(|X_1 - E(X_1)|^2)}{n\lambda^2}.$$

*Proof.* Using

$$S_n^2 = \sum_{k=1}^n X_k^2 + \sum_{j \neq k} X_j X_k,$$

<sup>3</sup><http://individual.utoronto.ca/jordanbell/notes/L0.pdf>, p. 4, Theorem 5.

<sup>4</sup>Charalambos D. Aliprantis and Kim C. Border, *Infinite Dimensional Analysis: A Hitchhikers Guide*, third ed., p. 485, Theorem 13.46.

we have

$$\begin{aligned}
E\left(\left|\frac{S_n}{n} - E(X_1)\right|^2\right) &= E\left(\sum_{k=1}^n \frac{X_k^2}{n^2} + \sum_{j \neq k} \frac{X_j X_k}{n^2} - 2E(X_1) \frac{S_n}{n} + E(X_1)^2\right) \\
&= \sum_{k=1}^n \frac{E(X_k^2)}{n^2} + \sum_{j \neq k} \frac{E(X_j X_k)}{n^2} - 2E(X_1)^2 + E(X_1)^2 \\
&= \sum_{k=1}^n \frac{E(X_1^2)}{n^2} + \sum_{j \neq k} \frac{E(X_j)E(X_k)}{n^2} - E(X_1)^2 \\
&= \sum_{k=1}^n \frac{E(X_1^2)}{n^2} + \sum_{j \neq k} \frac{E(X_1)^2}{n^2} + E(X_1)^2 \\
&= \frac{E(X_1^2)}{n} + \frac{E(X_1)^2}{n}.
\end{aligned}$$

On the other hand,

$$E(|X_1 - E(X_1)|^2) = E(X_1^2 - 2E(X_1)X_1 + E(X_1)^2) = E(X_1^2) + E(X_1)^2.$$

So

$$E\left(\left|\frac{S_n}{n} - E(X_1)\right|^2\right) = \frac{1}{n}E(|X_1 - E(X_1)|^2).$$

Using this and Chebyshev's inequality,

$$\begin{aligned}
P\left(\left|\frac{S_n}{n} - E(X_1)\right| \geq \lambda\right) &\leq \frac{1}{\lambda^2} E\left(\left|\frac{S_n}{n} - E(X_1)\right|^2\right) \\
&= \frac{1}{n\lambda^2} E(|X_1 - E(X_1)|^2),
\end{aligned}$$

proving the claim.  $\square$

We now prove the **weak law of large numbers**, which states that if  $X_n$  are independent identically distributed  $L^1$  random variables each with mean 0, then  $\frac{S_n}{n}$  converges in probability to  $E(X_1)$ .  $Z_n = X_n - E(X_n) = X_n - E(X_1)$  are independent and identically distributed  $L^1$  random variables with mean 0, and if  $Z_n$  converges in probability to 0 then  $X_n = Z_n + E(X_1)$  converges in probability to  $E(X_1)$ , showing that it suffices to prove the theorem when  $E(X_1) = 0$ .<sup>5</sup>

**Theorem 2** (Weak law of large numbers). *Suppose that  $E(X_1) = 0$ . For each  $\epsilon > 0$ ,*

$$P\left(\left|\frac{S_n}{n}\right| \geq \epsilon\right) \rightarrow 0, \quad n \rightarrow \infty.$$

<sup>5</sup>Allan Gut, *Probability: A Graduate Course*, second ed., p. 270, Theorem 6.3.1, and p. 121, Theorem 3.1.5.

*Proof.* For  $n \geq 1$  and  $1 \leq k \leq n$ , let

$$Y_{n,k} = 1_{\{|X_k| \leq n\epsilon^3\}} X_k = f \circ X_k,$$

where  $f(x) = 1_{[-n\epsilon^3, n\epsilon^3]}(x) \cdot x$ . Because  $X_1, \dots, X_k$  are independent and identically distributed, so are  $Y_{n,1}, \dots, Y_{n,n}$ .<sup>6</sup> Moreover,  $E(Y_{n,k}^2) \leq (n\epsilon^3)^2$ , so each  $Y_{n,k}$  belongs to  $L^2$ . Let

$$T_n = \sum_{k=1}^n Y_{n,k},$$

for which  $E(T_n) = nE(Y_{n,1})$ .

If  $\omega \in \bigcap_{k=1}^n \{|X_k| \leq n\epsilon^3\}$ , then

$$T_n(\omega) = \sum_{k=1}^n Y_{n,k}(\omega) = \sum_{k=1}^n X_k(\omega) = S_n(\omega),$$

and using this and Lemma 1, for  $t > 0$  we have

$$\begin{aligned} P(|S_n - E(T_n)| \geq t) &= P\left(\{|S_n - E(T_n)| \geq t\} \cap \bigcap_{k=1}^n \{|X_k| \leq n\epsilon^3\}\right) \\ &\quad + P\left(\{|S_n - E(T_n)| \geq t\} \cap \bigcup_{k=1}^n \{|X_k| > n\epsilon^3\}\right) \\ &\leq P(|T_n - E(T_n)| \geq t) + P\left(\bigcup_{k=1}^n \{|X_k| > n\epsilon^3\}\right) \\ &= P\left(\left|\frac{T_n}{n} - E(Y_{n,1})\right| \geq \frac{t}{n}\right) + P\left(\bigcup_{k=1}^n \{|X_k| > n\epsilon^3\}\right) \\ &\leq \frac{E(|Y_{n,1} - E(Y_{n,1})|^2)}{n \cdot \left(\frac{t}{n}\right)^2} + \sum_{k=1}^n P(|X_k| > n\epsilon^3) \\ &= \frac{n}{t^2} \cdot (E(Y_{n,1}^2) - E(Y_{n,1})^2) + nP(|X_1| > n\epsilon^3) \\ &\leq \frac{n}{t^2} \cdot E(Y_{n,1}^2) + nP(|X_1| > n\epsilon^3). \end{aligned}$$

For  $t = n\epsilon$  this is

$$\begin{aligned} P(|S_n - E(T_n)| \geq n\epsilon) &\leq \frac{1}{n\epsilon^2} \cdot E(Y_{n,1}^2) + nP(|X_1| > n\epsilon^3) \\ &= \frac{1}{n\epsilon^2} \cdot E(1_{\{|X_1| \leq n\epsilon^3\}} X_1^2) + nP(|X_1| > n\epsilon^3) \\ &\leq \frac{1}{n\epsilon^2} \cdot E(1_{\{|X_1| \leq n\epsilon^3\}} n\epsilon^3 |X_1|) + nP(|X_1| > n\epsilon^3) \\ &\leq \epsilon E(|X_1|) + nP(|X_1| > n\epsilon^3). \end{aligned}$$

<sup>6</sup>Gerald B. Folland, *Real Analysis: Modern Techniques and Their Applications*, second ed., p. 316, Proposition 10.2.

Now,

$$n\epsilon^3 P(|X_1| > n\epsilon^3) = n\epsilon^3 \int_{(n\epsilon^3, \infty)} d(X_{1*}P)(y) \leq \int_{(n\epsilon^3, \infty)} y d(X_{1*}P)(y),$$

which tends to 0 as  $n \rightarrow \infty$  because

$$\int_{\mathbb{R}} |y| d(X_{1*}P)(y) = E(|X_1|) < \infty.$$

Therefore,

$$\limsup_n P(|S_n - E(T_n)| \geq n\epsilon) \leq \epsilon E(|X_1|),$$

that is,

$$\limsup_n P\left(\left|\frac{S_n - E(T_n)}{n}\right| \geq \epsilon\right) \leq \epsilon E(|X_1|),$$

which implies that  $\frac{S_n - E(T_n)}{n}$  converges in probability to 0.

Because  $E(X_1) = 0$ ,

$$E(1_{\{|X_1| \leq n\epsilon^3\}} X_1) + E(1_{\{|X_1| > n\epsilon^3\}} X_1) = 0,$$

and hence

$$|E(T_n)| = |nE(Y_{n,1})| = n|E(1_{\{|X_1| \leq n\epsilon^3\}} X_1)| = n|E(1_{\{|X_1| > n\epsilon^3\}} X_1)|,$$

thus

$$\frac{|E(T_n)|}{n} = |E(1_{\{|X_1| > n\epsilon^3\}} X_1)| \leq E(1_{\{|X_1| > n\epsilon^3\}} |X_1|),$$

which tends to 0 as  $n \rightarrow \infty$ , because  $E(|X_1|) < \infty$ . Thus  $\frac{E(T_n)}{n}$  converges in probability to 0, and therefore  $\frac{S_n}{n}$  converges in probability to 0, completing the proof.  $\square$

**Lemma 3** (Bienaymé's formula). *If  $X_n$ ,  $n \geq 1$ , are  $L^2$  random variables that are pairwise uncorrelated, then*

$$\text{Var}\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n \text{Var}(X_k).$$

*Proof.* Let  $Y_k = X_k - E(X_k)$ . Using that the  $X_k$  are pairwise uncorrelated, for  $j \neq k$  we have

$$\begin{aligned} E(Y_j Y_k) &= E(X_j X_k - X_j E(X_k) - X_k E(X_j) + E(X_j) E(X_k)) \\ &= E(X_j) E(X_k) - E(X_j) E(X_k) - E(X_k) E(X_j) + E(X_j) E(X_k) \\ &= 0, \end{aligned}$$

showing that the  $Y_k$  are pairwise uncorrelated. Then, because  $E(S_n) = \sum_{k=1}^n E(X_k)$ ,

$$\begin{aligned}
\text{Var}(S_n) &= E\left(\left(S_n - \sum_{k=1}^n E(X_k)\right)^2\right) \\
&= E\left(\left(\sum_{k=1}^n Y_k\right)^2\right) \\
&= E\left(\sum_{k=1}^n Y_k^2 + \sum_{j \neq k} Y_j Y_k\right) \\
&= \sum_{k=1}^n E(Y_k^2) + \sum_{j \neq k} E(Y_j)E(Y_k) \\
&= \sum_{k=1}^n E(Y_k^2),
\end{aligned}$$

and as  $E(Y_k^2) = \text{Var}(X_k)$ , we have

$$\text{Var}(S_n) = \sum_{k=1}^n \text{Var}(X_k).$$

□

We now prove a weak law of large numbers, that is sometimes attributed to Markov, that neither supposes that the random variables are independent nor supposes that they are identically distributed. We remind ourselves that if a sequence  $Y_n$  of  $L^2$  random variables converges to 0 in  $L^2$  then it converges in probability to 0; this is proved using Chebyshev's inequality.

**Theorem 4** (Markov's weak law of large numbers). *If  $X_n$ ,  $n \geq 1$ , are  $L^2$  random variables that are pairwise uncorrelated and which satisfy*

$$\frac{1}{n^2} \sum_{k=1}^n \text{Var}(X_k) \rightarrow 0,$$

*then  $\frac{1}{n} \sum_{k=1}^n (X_k - E(X_k))$  converges to 0 in  $L^2$  and thus in probability.*

*Proof.* Because  $E(\sum_{k=1}^n X_k) = \sum_{k=1}^n E(X_k)$  and  $\text{Var}(X) = E((X - E(X))^2)$ ,

using Bienaymé's formula we get

$$\begin{aligned} E \left( \left( \frac{1}{n} \sum_{k=1}^n (X_k - E(X_k)) \right)^2 \right) &= \frac{1}{n^2} E \left( \left( \sum_{k=1}^n X_k - \sum_{k=1}^n E(X_k) \right)^2 \right) \\ &= \frac{1}{n^2} \text{Var} \left( \sum_{k=1}^n X_k \right) \\ &= \frac{1}{n^2} \sum_{k=1}^n \text{Var}(X_k). \end{aligned}$$

Thus

$$E \left( \left( \frac{1}{n} \sum_{k=1}^n (X_k - E(X_k)) \right)^2 \right) = \frac{1}{n^2} \sum_{k=1}^n \text{Var}(X_k) \rightarrow 0$$

as  $n \rightarrow \infty$ , namely,  $\frac{1}{n} \sum_{k=1}^n (X_k - E(X_k))$  converges to 0 in  $L^2$  as  $n \rightarrow \infty$ , proving the claim.  $\square$

### 3 Ergodic theory

We here assemble machinery and results that we will use to prove the strong law of large numbers. For a probability space  $(\Omega, \mathcal{F}, P)$ , a function  $T : \Omega \rightarrow \Omega$  is said to be a **measure-preserving transformation** if (i) it is measurable, and (ii)  $T_*P = P$ . To say that  $T_*P = P$  means that for any  $A \in \mathcal{F}$ ,  $(T_*P)(A) = P(A)$ , i.e.  $P(T^{-1}(A)) = P(A)$ .

A collection  $\mathcal{A}$  of subsets of  $\Omega$  is called an **algebra of sets** if (i)  $\emptyset \in \mathcal{A}$ , (ii)  $\Omega \in \mathcal{A}$ , (iii) if  $A, B \in \mathcal{A}$  then  $A \cup B, A \setminus B \in \mathcal{A}$ . If  $\mathcal{G}$  is a nonempty collection of subsets of  $\Omega$ , it is a fact that there is a unique algebra of sets  $A(\mathcal{G})$  that (i) contains  $\mathcal{G}$  and (ii) is contained in any algebra of sets that contains  $\mathcal{G}$ . We call  $A(\mathcal{G})$  the **algebra of sets generated by  $\mathcal{G}$** .

A collection  $\mathcal{S}$  of subsets of  $\Omega$  is called a **semialgebra of sets** if (i)  $\emptyset \in \mathcal{S}$ , (ii)  $\Omega \in \mathcal{S}$ , (iii) if  $A, B \in \mathcal{S}$  then  $A \cap B \in \mathcal{S}$ , (iv) if  $A, B \in \mathcal{S}$  then there are pairwise disjoint  $E_1, \dots, E_n \in \mathcal{S}$  such that

$$A \setminus B = \bigcup_{k=1}^n E_k.$$

It is a fact that  $A(\mathcal{S})$  is equal to the collection of all unions of finitely many pairwise disjoint elements of  $\mathcal{S}$ .<sup>7</sup>

A nonempty collection  $\mathcal{M}$  of subsets of  $\Omega$  is called a **monotone class** if whenever  $A_n \in \mathcal{M}$ ,  $A_n \subset A_{n+1}$  (an **increasing sequence of sets**), implies that  $\bigcup_n A_n \in \mathcal{M}$  and  $A_n \in \mathcal{M}$ ,  $A_{n+1} \subset A_n$  (a **decreasing sequence of sets**), implies that  $\bigcap_n A_n \in \mathcal{M}$ . In other words, a monotone class is a nonempty

<sup>7</sup>V. I. Bogachev, *Measure Theory*, volume I, p. 8, Lemma 1.2.14.

collection of subsets of  $\Omega$  such that if  $A_n$  is a monotone sequence in  $\mathcal{M}$  then  $\lim_{n \rightarrow \infty} A_n \in \mathcal{M}$ . If  $\mathcal{G}$  is a nonempty collection of subsets of  $\Omega$ , it is a fact that there is a unique monotone class  $M(\mathcal{G})$  that (i) contains  $\mathcal{G}$  and (ii) is contained in any monotone class that contains  $\mathcal{G}$ . We call  $M(\mathcal{G})$  the **monotone class generated by  $\mathcal{G}$** . The **monotone class theorem** states that if  $\mathcal{A}$  is an algebra of sets then  $\sigma(\mathcal{A}) = M(\mathcal{A})$ .<sup>8</sup>

The following lemma gives conditions under which we can establish that a function is measure-preserving.<sup>9</sup>

**Lemma 5.** *Let  $T : \Omega \rightarrow \Omega$  be a function and suppose that  $\mathcal{S}$  is a semialgebra of sets for which  $\mathcal{F}$  is the  $\sigma$ -algebra generated by  $\mathcal{S}$ . If (i)  $T^{-1}(A) \in \mathcal{F}$  for each  $A \in \mathcal{S}$  and (ii)  $P(T^{-1}(A)) = P(A)$  for each  $A \in \mathcal{S}$ , then  $T$  is measure-preserving.*

*Proof.* Let

$$\mathcal{C} = \{A \in \mathcal{F} : T^{-1}(A) \in \mathcal{F}, P(T^{-1}(A)) = P(A)\};$$

we wish to prove that  $\mathcal{C} = \mathcal{F}$ .

If  $A_n \in \mathcal{C}$  is an increasing sequence, let  $A = \bigcup_{n=1}^{\infty} A_n$ . Then, as  $T^{-1}(A_n) \in \mathcal{F}$  for each  $n$ ,

$$T^{-1}(A) = \bigcup_{n=1}^{\infty} T^{-1}(A_n) \in \mathcal{F}$$

and as (i)  $T^{-1}(A_n)$  is an increasing sequence, (ii)  $P$  is continuous from below,<sup>10</sup> and (iii)  $P(T^{-1}(A_n)) = P(A_n)$ ,

$$\begin{aligned} P(T^{-1}(A)) &= P\left(\bigcup_{n=1}^{\infty} T^{-1}(A_n)\right) \\ &= \lim_{n \rightarrow \infty} P(T^{-1}(A_n)) \\ &= \lim_{n \rightarrow \infty} P(A_n) \\ &= P\left(\bigcup_{n=1}^{\infty} A_n\right) \\ &= P(A), \end{aligned}$$

and hence  $A \in \mathcal{C}$ . If  $A_n \in \mathcal{C}$  is a decreasing sequence, let  $A = \bigcap_{n=1}^{\infty} A_n$ . Because  $T^{-1}(A_n) \in \mathcal{F}$ ,

$$T^{-1}(A) = \bigcap_{n=1}^{\infty} T^{-1}(A_n) \in \mathcal{F},$$

<sup>8</sup>V. I. Bogachev, *Measure Theory*, volume I, p. 33, Theorem 1.9.3.

<sup>9</sup>Peter Walters, *An Introduction to Ergodic Theory*, p. 20, Theorem 1.1.

<sup>10</sup>V. I. Bogachev, *Measure Theory*, volume I, p. 9, Proposition 1.3.3.



and as (i)  $T^{-1}(A_n)$  is a decreasing sequence, (ii)  $P$  is continuous from above, and (iii)  $P(T^{-1}(A_n)) = P(A_n)$ ,

$$\begin{aligned} P(T^{-1}(A)) &= P\left(\bigcap_{n=1}^{\infty} T^{-1}(A_n)\right) \\ &= \lim_{n \rightarrow \infty} P(T^{-1}(A_n)) \\ &= \lim_{n \rightarrow \infty} P(A_n) \\ &= P\left(\bigcap_{n=1}^{\infty} A_n\right) \\ &= P(A), \end{aligned}$$

and hence  $A \in \mathcal{C}$ . Therefore,  $\mathcal{C}$  is a monotone class.

$\mathcal{S} \subset \mathcal{C}$ . If  $A \in A(\mathcal{S})$ , then there are pairwise disjoint  $A_1, \dots, A_n \in \mathcal{S}$  with  $A = \bigcup_{k=1}^n A_k$ . As  $T^{-1}(A_k) \in \mathcal{F}$ ,

$$T^{-1}(A) = \bigcup_{k=1}^n T^{-1}(A_k) \in \mathcal{F}.$$

As the  $A_k$  are pairwise disjoint, so are the sets  $T^{-1}(A_k)$ , so, because  $P(T^{-1}(A_k)) = P(A_k)$ ,

$$\begin{aligned} P(T^{-1}(A)) &= P\left(\bigcup_{k=1}^n T^{-1}(A_k)\right) \\ &= \sum_{k=1}^n P(T^{-1}(A_k)) \\ &= \sum_{k=1}^n P(A_k) \\ &= P\left(\bigcup_{k=1}^n A_k\right) \\ &= P(A). \end{aligned}$$

Therefore  $A \in \mathcal{C}$ . This shows that  $A(\mathcal{S}) \subset \mathcal{C}$ .

The monotone class theorem tells us  $\sigma(A(\mathcal{S})) = M(A(\mathcal{S}))$ . On the one hand,  $\mathcal{F} = \sigma(\mathcal{S})$  and hence  $\mathcal{F} = \sigma(A(\mathcal{S}))$ . On the other hand,  $A(\mathcal{S}) \subset \mathcal{C}$  and the fact that  $\mathcal{C}$  is a monotone class yield

$$M(A(\mathcal{S})) \subset M(\mathcal{C}) = \mathcal{C}.$$

Therefore

$$\mathcal{F} \subset \mathcal{C}.$$

Of course  $\mathcal{C} \subset \mathcal{F}$ , so  $\mathcal{C} = \mathcal{F}$ , proving the claim.  $\square$

Let  $(\Omega, \mathcal{F}, P)$  be a probability space. A measure-preserving transformation  $T : \Omega \rightarrow \Omega$  is called **ergodic** if  $A \in \mathcal{F}$  and  $T^{-1}(A) = A$  implies that  $P(A) = 0$  or  $P(A) = 1$ . It is proved<sup>11</sup> using the **Birkhoff ergodic theorem** that for a measure-preserving transformation  $T : \Omega \rightarrow \Omega$ ,  $T$  is ergodic if and only if for all  $A, B \in \mathcal{F}$ ,

$$\frac{1}{n} \sum_{k=0}^{n-1} P(T^{-k}(A) \cap B) \rightarrow P(A)P(B).$$

A measure-preserving transformation  $T : \Omega \rightarrow \Omega$  is called **weak-mixing** if for all  $A, B \in \mathcal{F}$ ,

$$\frac{1}{n} \sum_{k=0}^{n-1} |P(T^{-k}(A) \cap B) - P(A)P(B)| \rightarrow 0.$$

It is immediate that a weak-mixing transformation is ergodic.

A measure-preserving transformation  $T : \Omega \rightarrow \Omega$  is called **strong-mixing** if for all  $A, B \in \mathcal{F}$ ,

$$P(T^{-n}(A) \cap B) \rightarrow P(A)P(B).$$

If a sequence of real numbers  $a_n$  tends to 0, then

$$\frac{1}{n} \sum_{k=0}^{n-1} |a_k| \rightarrow 0,$$

and using this we check that a strong-mixing transformation is weak-mixing.

The following statement gives conditions under which a measure-preserving transformation is ergodic, weak-mixing, or strong-mixing.<sup>12</sup>

**Theorem 6.** *Let  $(\Omega, \mathcal{F}, P)$  be a probability space, let  $\mathcal{S}$  be a semialgebra that generates  $\mathcal{F}$ , and let  $T : \Omega \rightarrow \Omega$  be a measure-preserving transformation.*

1.  *$T$  is ergodic if and only if for all  $A, B \in \mathcal{S}$ ,*

$$\frac{1}{n} \sum_{k=0}^{n-1} P(T^{-k}(A) \cap B) \rightarrow P(A)P(B).$$

2.  *$T$  is weak-mixing if and only if for all  $A, B \in \mathcal{S}$ ,*

$$\frac{1}{n} \sum_{k=0}^{n-1} |P(T^{-k}(A) \cap B) - P(A)P(B)| \rightarrow 0.$$

3.  *$T$  is strong-mixing if and only if for all  $A, B \in \mathcal{S}$ ,*

$$P(T^{-n}(A) \cap B) \rightarrow P(A)P(B).$$

<sup>11</sup>Peter Walters, *An Introduction to Ergodic Theory*, p. 37, Corollary 1.14.2.

<sup>12</sup>Peter Walters, *An Introduction to Ergodic Theory*, p. 41, Theorem 1.17.

## 4 The strong law of large numbers

Let  $\mu$  be a Borel probability measure on  $\mathbb{R}$  with **finite first moment**:

$$\int_{\mathbb{R}} |x| dm(x) < \infty.$$

We shall specify when we use the hypothesis that  $\mu$  has finite first moment; until we say so, what we say merely supposes that it is a Borel probability measure on  $\mathbb{R}$ .

For  $n \geq 0$ , let  $\Omega_n = \mathbb{R}$ , let  $\mathcal{B}_n = \mathcal{B}_{\mathbb{R}}$ , the Borel  $\sigma$ -algebra of  $\mathbb{R}$ , and let  $\mu_n = \mu$ , for which  $(\Omega_n, \mathcal{B}_n, \mu_n)$  is a probability space. Let  $\Omega = \prod_{n=0}^{\infty} \Omega_n$ . A **cylinder set** is a subset of  $\Omega$  of the form

$$\prod_{n=0}^{\infty} A_n,$$

where  $A_n \in \mathcal{B}_n$  for each  $n$  and where  $\{n \geq 0 : A_n \neq \Omega_n\}$  is finite. We denote by  $\mathcal{C}$  the collection of all cylinder sets. It is a fact that  $\mathcal{C}$  is a semialgebra of sets.<sup>13</sup>

The **product  $\sigma$ -algebra** is  $\mathcal{F} = \sigma(\mathcal{C})$ , the  $\sigma$ -algebra generated by the collection of all cylinder sets. The **product measure**,<sup>14</sup> which we denote by  $P$ , is the unique probability measure on  $\mathcal{F}$  such that for any cylinder set  $\prod_{n=0}^{\infty} A_n$ ,

$$P\left(\prod_{n=0}^{\infty} A_n\right) = \prod_{n=0}^{\infty} \mu_n(A_n).$$

Let  $\pi_n : \Omega \rightarrow \Omega_n$ ,  $n \geq 0$ , be the projection map. Define  $\tau : \Omega \rightarrow \Omega$  by

$$\tau(\omega_0, \omega_1, \dots) = (\omega_1, \omega_2, \dots),$$

the **left shift map**. In other words, for  $n \geq 0$ ,  $\tau$  satisfies  $\pi_n \circ \tau = \pi_{n+1}$ , and so  $\pi_n = \pi_0 \circ (\tau^n)$ .

**Lemma 7.**  $\tau$  is measure-preserving.

*Proof.* For  $A = \prod_{n=0}^{\infty} A_n \in \mathcal{C}$ ,

$$\tau^{-1}(A) = \prod_{n=0}^{\infty} B_n = B,$$

where  $B_0 = \Omega_0$  and for  $n \geq 1$ ,  $B_n = A_{n-1}$ .  $B$  is a cylinder set so a fortiori belongs to  $\mathcal{F}$ , and

$$P(B) = \prod_{n=0}^{\infty} \mu_n(B_n) = \mu_0(\Omega_0) \cdot \prod_{n=1}^{\infty} \mu_n(B_n) = \prod_{n=1}^{\infty} \mu_n(A_{n-1}) = \prod_{n=1}^{\infty} \mu_{n-1}(A_{n-1}),$$

<sup>13</sup>S. J. Taylor, *Introduction to Measure Theory and Integration*, p. 136, §6.1; <http://individual.utoronto.ca/jordanbell/notes/productmeasure.pdf>

<sup>14</sup><http://individual.utoronto.ca/jordanbell/notes/productmeasure.pdf>

so

$$P(\tau^{-1}(A)) = \prod_{n=0}^{\infty} \mu_n(A_n) = P(A).$$

Therefore by Lemma 5, because  $\mathcal{C}$  is a semialgebra that generates  $\mathcal{F}$ , it follows that  $\tau$  is measure-preserving.  $\square$

**Lemma 8.**  $\tau$  is strong-mixing.

*Proof.* Let  $A = \prod_{n=0}^{\infty} A_n$  and  $B = \prod_{n=0}^{\infty} B_n$  be cylinder sets. For  $n \geq 0$ ,

$$\tau^{-n}(A) = \prod_{m=0}^{\infty} C_m,$$

where  $C_m = \Omega_m$  for  $0 \leq m \leq n-1$  and  $C_m = A_{m-n}$  for  $m \geq n$ . Because  $A$  and  $B$  are cylinder sets, there is some  $N$  such that when  $m \geq N$ ,  $A_m = \Omega_m$  and  $B_m = \Omega_m$ . Thus for  $n \geq N$ ,

$$\begin{aligned} \tau^{-n}(A) \cap B &= \prod_{m=0}^{\infty} C_m \cap \prod_{m=0}^{\infty} B_m \\ &= \prod_{m=0}^{\infty} (C_m \cap B_m) \\ &= \prod_{m=0}^{n-1} (C_m \cap B_m) \times \prod_{m=n}^{\infty} (C_m \cap B_m) \\ &= \prod_{m=0}^{n-1} (\Omega_m \cap B_m) \times \prod_{m=n}^{\infty} (A_{m-n} \cap \Omega_m) \\ &= \prod_{m=0}^{n-1} B_m \times \prod_{m=n}^{\infty} A_{m-n}. \end{aligned}$$

Hence

$$\begin{aligned} P(\tau^{-n}(A) \cap B) &= \prod_{m=0}^{n-1} \mu_m(B_m) \cdot \prod_{m=n}^{\infty} \mu_m(A_{m-n}) \\ &= \prod_{m=0}^{n-1} \mu_m(B_m) \cdot \prod_{m=0}^{\infty} \mu_{m+n}(A_m) \\ &= \prod_{m=0}^{n-1} \mu_m(B_m) \cdot \prod_{m=0}^{\infty} \mu_m(A_m) \\ &= P(B) \cdot P(A). \end{aligned}$$

That is, there is some  $N$  such that when  $n \geq N$ ,

$$P(\tau^{-n}(A) \cap B) = P(A)P(B),$$

and so a fortiori,

$$\lim_{n \rightarrow \infty} P(\tau^{-n}(A) \cap B) = P(A)P(B).$$

Therefore, because the cylinder sets generate the  $\sigma$ -algebra  $\mathcal{F}$ , by Theorem 3 we get that  $\tau$  is strong-mixing.  $\square$

We now use the hypothesis that  $\mu$  has finite first moment.<sup>15</sup>

**Lemma 9.**  $\pi_0 \in L^1(P)$ .

*Proof.*  $T = \pi_0 : \Omega \rightarrow \Omega_0$  is measurable, and  $T_*P = \mu_0$ . The statement that  $\mu = \mu_0$  has finite first moment means that  $f : \Omega_0 \rightarrow \mathbb{R}$  defined by  $f(x) = |x|$  belongs to  $L^1(\mu_0)$ . Therefore by the change of variables theorem (2), we have  $f \circ T \in L^1(P)$ .  $\square$

**Lemma 10.** For almost all  $\omega \in \Omega$ ,

$$\frac{1}{n} \sum_{k=0}^{n-1} \pi_k(\omega) \rightarrow \int_{\mathbb{R}} t d\mu(t).$$

*Proof.* Because  $\tau : \Omega \rightarrow \Omega$  is ergodic and  $\pi_0 \in L^1(P)$ , the **Birkhoff ergodic theorem**<sup>16</sup> tells us that for almost all  $\omega \in \Omega$ ,

$$\frac{1}{n} \sum_{k=0}^{n-1} \pi_0(\tau^k(\omega)) \rightarrow \int_{\Omega} \pi_0(\omega) dP(\omega),$$

i.e.,

$$\frac{1}{n} \sum_{k=0}^{n-1} \pi_k(\omega) \rightarrow \int_{\Omega_0} \omega_0 d\mu_0(\omega_0),$$

proving the claim.  $\square$

We will use Lemma 10 to prove the strong law of large numbers. First we prove two lemmas about **joint distributions**.<sup>17</sup>

**Lemma 11.** If  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  and  $Y_1, \dots, Y_n : \Omega' \rightarrow \mathbb{R}$  are random variables with the same joint distribution and for each  $1 \leq k \leq n$ ,  $\Phi_k : \mathbb{R}^n \rightarrow \mathbb{R}$  is a continuous function, then for  $U_k = \Phi_k(X_1, \dots, X_n)$  and  $V_k = \Phi_k(Y_1, \dots, Y_n)$ , the random variables  $U_1, \dots, U_n$  and  $V_1, \dots, V_n$  have the same joint distribution.

<sup>15</sup>Elias M. Stein and Rami Shakarchi, *Functional Analysis*, Princeton Lectures in Analysis, volume IV, p. 208, chapter 5, §2.1.

<sup>16</sup>Peter Walters, *An Introduction to Ergodic Theory*, p. 34, Theorem 1.14.

<sup>17</sup>Elias M. Stein and Rami Shakarchi, *Functional Analysis*, Princeton Lectures in Analysis, volume IV, p. 208, Lemma 2.2, Lemma 2.3.

*Proof.* Write  $X = (X_1, \dots, X_n)$ ,  $Y = (Y_1, \dots, Y_n)$ ,  $U = (U_1, \dots, U_n)$ , and  $V = (V_1, \dots, V_n)$ , which are Borel measurable. Let  $\Phi = (\Phi_1, \dots, \Phi_n)$ , which is continuous  $\mathbb{R}^n \rightarrow \mathbb{R}^n$ , and for which

$$U = \Phi(X), \quad V = \Phi(Y).$$

To show that  $U_1, \dots, U_n$  and  $V_1, \dots, V_n$  have the same joint distribution means to show that  $U_*P = V_*P'$ . Let  $A \in \mathcal{B}_{\mathbb{R}^n}$ , for which  $\Phi^{-1}(A) \in \mathcal{B}_{\mathbb{R}^n}$  and

$$\begin{aligned} (U_*P)(A) &= P(U^{-1}(A)) \\ &= P(X^{-1}(\Phi^{-1}(A))) \\ &= P'(Y^{-1}(\Phi^{-1}(A))) \\ &= P'(V^{-1}(A)) \\ &= (V_*P')(A), \end{aligned}$$

where, because  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  have the same joint distribution,

$$P(X^{-1}(\Phi^{-1}(A))) = (X_*P)(\Phi^{-1}(A)) = (Y_*P')(\Phi^{-1}(A)) = P'(Y^{-1}(\Phi^{-1}(A))).$$

□

**Lemma 12.** *If sequences of random variables  $X_n : (\Omega, \mathcal{F}, P) \rightarrow \mathbb{R}$  and  $Y_n : (\Omega', \mathcal{F}', P') \rightarrow \mathbb{R}$ ,  $n \geq 1$ , have the same **finite-dimensional distributions** and there is some  $a \in \mathbb{R}$  such that  $X_n$  converges to  $a$  almost surely, then  $Y_n$  converges to  $a$  almost surely.*

*Proof.* That the sequences  $X_n$  and  $Y_n$  have the same finite-dimensional distributions means that for each  $n \geq 1$ ,

$$(X_1, \dots, X_n)_*P = (Y_1, \dots, Y_n)_*P',$$

i.e., for each  $n \geq 1$  and for each  $A \in \mathcal{B}_{\mathbb{R}^n}$ ,

$$P((X_1, \dots, X_n) \in A) = P'((Y_1, \dots, Y_n) \in A).$$

Define

$$\begin{aligned} E_{k,N,n} &= \left\{ \omega \in \Omega : |X_n(\omega) - a| \leq \frac{1}{k} \right\} \\ F_{k,N,n} &= \left\{ \omega \in \Omega' : |Y_n(\omega) - a| \leq \frac{1}{k} \right\}. \end{aligned}$$

Then define

$$\begin{aligned} E &= \bigcap_{k=1}^{\infty} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} E_{k,N,n} = \bigcap_{k=1}^{\infty} \bigcup_{N=1}^{\infty} E_{k,N} = \bigcap_{k=1}^{\infty} E_k \\ F &= \bigcap_{k=1}^{\infty} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} F_{k,N,n} = \bigcap_{k=1}^{\infty} \bigcup_{N=1}^{\infty} F_{k,N} = \bigcap_{k=1}^{\infty} F_k, \end{aligned}$$

and

$$G_{k,N,n} = \bigcap_{m=N}^n E_{k,N,m}$$

$$H_{k,N,n} = \bigcap_{m=N}^n F_{k,N,m}.$$

Because  $(X_1, \dots, X_n)$  and  $(Y_1, \dots, Y_n)$  have the same distribution,

$$P(G_{k,N,n}) = P'(H_{k,N,n}).$$

But

$$E_{k,N} = \bigcap_{n=N}^{\infty} E_{k,N,n} = \bigcap_{n=N}^{\infty} G_{k,N,n} = \lim_{n \rightarrow \infty} G_{k,N,n}$$

and

$$F_{k,N} = \lim_{n \rightarrow \infty} H_{k,N,n},$$

so

$$P(E_{k,N}) = P'(F_{k,N}).$$

Then

$$P(E_k) = \lim_{N \rightarrow \infty} P(E_{k,N}) = \lim_{N \rightarrow \infty} P'(F_{k,N}) = P'(F_k).$$

Then

$$P(E) = \lim_{k \rightarrow \infty} P(E_k) = \lim_{k \rightarrow \infty} P'(F_k) = P'(F).$$

That the sequence  $X_n$  converges almost surely means that  $P(E) = 1$ , and therefore  $P'(F) = 1$ , i.e. the sequence  $Y_n$  converges almost surely.  $\square$

We now use what we have established to prove the **strong law of large numbers**.<sup>18</sup> (We write  $(\Omega', \mathcal{F}', P')$  for a probability space because  $(\Omega, \mathcal{F}, P)$  denotes the product probability space constructed already.)

**Theorem 13.** *If  $X_n : (\Omega', \mathcal{F}', P') \rightarrow \mathbb{R}$ ,  $n \geq 1$ , are independent identically distributed  $L^1$  random variables, then for almost all  $\omega \in \Omega'$ ,*

$$\frac{1}{n} \sum_{k=1}^n X_k(\omega) \rightarrow E(X_1).$$

*Proof.* For  $n \geq 0$ , let  $Y_n = X_{n+1}$ ; we do this to make the index set the same as for  $\Omega_n$ . Let  $\mu_n = Y_{n*}P'$  and set  $\mu = \mu_0$ . Because the  $Y_n$  are identically distributed,  $\mu_n = \mu_0$  for each  $n$ .

Because  $Y_0$  is  $L^1$ ,  $\mu$  has finite first moment. As  $\mu_0 = Y_{0*}P'$ , applying the change of variables theorem (2),

$$\int_{\mathbb{R}} td\mu(t) = \int_{\mathbb{R}} td(Y_{0*}P')(t) = \int_{\Omega'} Y_0(\omega)dP'(\omega) = E(Y_0).$$

<sup>18</sup>Elias M. Stein and Rami Shakarchi, *Functional Analysis*, Princeton Lectures in Analysis, volume IV, p. 206, Theorem 2.1.

With this, Lemma 10 says that for almost all  $\omega \in \Omega$ ,

$$\frac{1}{n} \sum_{k=0}^{n-1} \pi_k(\omega) \rightarrow E(Y_0). \quad (3)$$

For each  $n$ ,

$$(\pi_0, \dots, \pi_n)_* P = \mu_0 \times \dots \times \mu_n,$$

and because the  $Y_n$  are independent,

$$(Y_0, \dots, Y_n)_* P' = Y_0_* P' \times \dots \times Y_n_* P' = \mu_0 \times \dots \times \mu_n,$$

so the sequences  $\pi_n$  and  $Y_n$  have the same finite-dimensional distributions. For  $n \geq 1$ , define  $\Phi_n : \mathbb{R}^n \rightarrow \mathbb{R}$  by

$$\Phi_n(t_1, \dots, t_n) = \frac{1}{n} \sum_{k=1}^n t_k.$$

Lemma 11 then tells us that the sequences

$$U_n = \Phi_n(Y_0, \dots, Y_{n-1}) = \frac{1}{n} \sum_{k=0}^{n-1} Y_k$$

and

$$V_n = \Phi_n(\pi_0, \dots, \pi_{n-1}) = \frac{1}{n} \sum_{k=0}^{n-1} \pi_k,$$

$n \geq 1$ , have the same finite-dimensional distributions.

Now, (3) says that  $V_n$  converges to  $E(Y_0)$  almost surely, and thus applying Lemma 12 we get that  $U_n$  converges to  $E(Y_0)$  almost surely. That is,

$$\frac{1}{n} \sum_{k=0}^{n-1} Y_k \rightarrow E(Y_0)$$

almost surely, and because  $Y_k = X_{k+1}$ ,

$$\frac{1}{n} \sum_{k=1}^n X_k \rightarrow E(X_1)$$

almost surely, completing the proof.  $\square$