

# Probability Notes 2024

Ruan Yuanlong

BUAA, BEIJING, CHINA

*Email address:* ruanyl@buaa.edu.cn



# Contents

1. 单调类定理	7
2. 集函数与测度	13
2.1. 集函数	13
2.2. 半环上非负集函数	17
2.3. 环上非负集函数	23
3. Carathéodory's 延拓	28
3.1. 外测度	28
3.2. 域上测度的延拓	35
3.3. 半环上测度的延拓	40
3.4. Approximating $\mu^* _{\mathcal{F}_\mu^*}$ by $\mu^* _{\sigma(\mathcal{S})}$	42

3.5.	Approximating $\mu _{\sigma(\mathcal{A})}$ by $\mu _{\mathcal{A}}$	44
3.6.	Completion of a measure space	44
4.	收敛	49
4.1.	可测函数的收敛	49
4.2.	随机变量的分布函数	57
4.3.	随机变量的收敛	59
5.	积分	64
5.1.	非负可测函数积分	64
5.2.	可测函数积分	68
5.3.	Change of variables	73
6.	$L_p$ 空间	73
6.1.	Inequalities	73
6.2.	Completeness	84
6.3.	$L_p$ and weak convergence	87
6.4.	Uniform integrability	93
6.5.	Summary of various convergences	100
7.	概率空间的积分	101
7.1.	Expected value	101

7.2.	Properties of expectation	106
7.3.	Lebesgue-Stieltjes and Riemann-Stieltjes integrals	108
7.4.	$L_p$ convergence and uniform integrability	112
8.	乘积测度空间	115
8.1.	Product $\sigma$ -field	115
8.2.	Product measure space	118
8.3.	Fubini's Theorem	123
8.4.	Applications	125
8.5.	Finite-dimensional product space	130
9.	独立性 Independence	132
9.1.	Independence of events and random variables	132
9.2.	Independence and expectation	143
9.3.	Sum of independent random variables	148
9.4.	Construction of independent sequence	161
10.	大数律 Law of large numbers	171
10.1.	$L_2$ weak law	171
10.2.	Weak law of large numbers	182
10.3.	Borel-Cantelli lemma and applications	194

10.4.	Strong law of large numbers	204
11.	中心极限定理 Central limit theorem	224
11.1.	Introduction	224
11.2.	From Poisson distribution to Stirling formula	225
11.3.	De Moivre-Laplace limit theorem	228
11.4.	Weak convergence	231
12.	Characteristic function	252
12.1.	A breif review of complex calculus	252
12.2.	The definition of Characteristic function	254
12.3.	The inversion formula and uniqueness	257
12.4.	Moments and derivatives	266
12.5.	Continuity theorem	273
13.	Central limit theorem – the proof	277

## 1. 单调类定理

Review:

- $\mathcal{A}$  is a field,  $\mathcal{M}$  is a monotone class. Then

$$\mathcal{A} \subset \mathcal{M} \implies \sigma(\mathcal{A}) \subset \mathcal{M}.$$

- $\mathcal{P}$  is a  $\pi$ -system,  $\mathcal{L}$  is a  $\lambda$ -system. Then

$$\mathcal{P} \subset \mathcal{L} \implies \sigma(\mathcal{P}) \subset \mathcal{L}.$$

- measurable spaces  $(E, \mathcal{F}_E), (F, \mathcal{F}_F), f : (E, \mathcal{F}_E) \mapsto (F, \mathcal{F}_F)$ .  
 $f$  is  $\mathcal{F}_E/\mathcal{F}_F$ -measurable if

$$\sigma(f) \triangleq f^{-1}(\mathcal{F}_F) \subset \mathcal{F}_E.$$

Call it  $\mathcal{F}_E$ -measurable if

$$(F, \mathcal{F}_F) = (\mathbb{R}, \mathcal{B}(\mathbb{R})).$$

- $f : (E, \mathcal{F}_E) \mapsto (F, \sigma(\mathcal{E}))$ ,  $f$  is  $\mathcal{F}_E/\sigma(\mathcal{E})$ -measurable if

$$f^{-1}(\mathcal{E}) \subset \mathcal{F}_E.$$

**Thm 1** ( $\pi$ - $\lambda$  theorem).  $\mathcal{P}$  is a  $\pi$ -system,  $\mathcal{L}$  is a  $\lambda$ -system. If  $\mathcal{P} \subset \mathcal{L}$ , then  $\sigma(\mathcal{P}) \subset \mathcal{L}$ .

**Def 1 (Simple function)**.  $i = 1, \dots, n$ ,  $A_i \in \mathcal{F}$  (pairwise) disjoint,  $c_i \in \mathbb{R}$ .  $f$  is (measurable) simple if  $f = \sum_{i=1}^n c_i 1_{A_i}$ .

**Alt.**  $i = 1, \dots, n$ ,  $A_i \in \mathcal{F}$ ,  $c_i \in \mathbb{R}$  non-zero distinct,  $f$  is simple if  $f = \sum_{i=1}^n c_i 1_{A_i}$ .

▷ 1.  $a, b \in \mathbb{R}$ ,  $g$  simple, then  $af + bg$  simple

**Thm 2 (Simple approximation)**. (1)  $f \geq 0$  measurable. There exist simple  $\{f_n\}$ ,  $0 \leq f_n \uparrow f$ , uniform if  $f$  is bounded.

(2)  $f$  measurable. There exist simple  $\{f_n\}$ ,  $f_n \rightarrow f$ , uniform if  $f$  is bounded.



PROOF. 1. Let

$$f_n = \frac{[2^n f]}{2^n} \wedge n = \sum_{i=0}^{n2^n-1} \frac{i}{2^n} 1_{\{i/2^n \leq f < (i+1)/2^n\}} + n 1_{\{f \geq n\}}.$$

Then

$$0 \leq f - f_n \leq \frac{1}{2^n} \text{ if } f < n; \quad f_n = n \leq f \text{ otherwise.}$$

2.  $f = f^+ - f^-$ . □

**Thm 3 (Doob).**  $f : (E, \mathcal{F}_E) \mapsto (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ ,  $g$  measurable  $(E, \mathcal{F}_E) \mapsto (F, \mathcal{F}_F)$ . If  $f$  is  $\sigma(g)$ -measurable, then  $f = h \circ g$  for some measurable  $h$ .

PROOF. 1.  $f = 1_A$ ,  $A = g^{-1}(B) \in \sigma(g)$ ,  $B \in \mathcal{F}_F$ . Then  $x \in A$  if and only if  $g(x) \in B$ , i.e.,

$$f = 1_A = 1_B \circ g.$$

**2.**  $f$  simple,  $f = \sum_{i=1}^n c_i 1_{A_i}$ ,  $c_i \in \mathbb{R}$ ,  $A_i \in \sigma(g)$  disjoint. Let  $A_i = g^{-1}(B_i)$ ,  $B_i \in \mathcal{F}_F$ , then

$$C_i = B_i \setminus \left( \bigcup_{j < i} B_j \right) \in \mathcal{F}_F \text{ disjoint}$$

and

$$f^{-1}(C_i) = A_i \setminus \left( \bigcup_{j < i} A_j \right) = A_i.$$

By step 1,

$$f = \sum_{i=1}^n c_i 1_{A_i} = \sum_{i=1}^n c_i 1_{C_i} \circ g = \left( \sum_{i=1}^n c_i 1_{C_i} \right) \circ g \triangleq h \circ g.$$

**3.**  $f \geq 0$  is  $\sigma(g)$ -measurable, there exist  $\sigma(g)$ -measurable simple  $f_n$  with  $0 \leq f_n \uparrow f$ . It follows  $f_n = h_n \circ g$  for some  $h_n$ ,

$$h \triangleq \sup_n h_n$$

is  $\sigma(g)$ -measurable,

$$f = \lim_n f_n = \sup_n (h_n \circ g) = \left( \sup_n h_n \right) \circ g = h \circ g.$$

4.  $f$  is  $\sigma(g)$ -measurable.  $f^+, f^-$  are  $\sigma(g)$ -measurable. Use **3**.  $\square$

**Thm 4.**  $\mathcal{A}$  is a  $\pi$ -system,  $\Omega \in \mathcal{A}$ ,  $\mathcal{H}$  is a collection of real-valued functions. Suppose

(1) If  $A \in \mathcal{A}$ , then  $1_A \in \mathcal{H}$

(2) If  $f, g \in \mathcal{H}$ ,  $c \in \mathbb{R}$ , then  $f + g, cg \in \mathcal{H}$

(3) If  $f_n \in \mathcal{H}$ ,  $0 \leq f_n \uparrow f$  with  $f$  bounded, then  $f \in \mathcal{H}$

Then

$$\{f : f \text{ bounded } \sigma(\mathcal{A})\text{-measurable}\} \subset \mathcal{H}$$

PROOF. The system of sets

$$\mathcal{G} = \{A : 1_A \in \mathcal{H}\}$$

is a  $\lambda$ -system and  $\mathcal{A} \subset \mathcal{G}$ . Hence

$$\sigma(\mathcal{A}) \subset \mathcal{G}.$$

(2) implies that  $\mathcal{H}$  contains all  $\sigma(\mathcal{A})$ -measurable simple functions, (3) implies that  $\mathcal{H}$  contains all bounded  $\sigma(\mathcal{A})$ -measurable functions.  $\square$

**Thm 5.**  $\mathcal{A}$  is a  $\pi$ -system,  $\Omega \in \mathcal{A}$ ,  $\mathcal{H}$  is a collection of real-valued functions. Suppose

(1) If  $A \in \mathcal{A}$ , then  $1_A \in \mathcal{H}$

(2) If  $f, g \in \mathcal{H}$ ,  $a, b \geq 0$ , then  $af + bg \in \mathcal{H}$

(3) If  $f, g \in \mathcal{H}$  are bounded,  $f \geq g$ , then  $f - g \in \mathcal{H}$

(4) If  $f_n \in \mathcal{H}$ ,  $0 \leq f_n \uparrow f$ , then  $f \in \mathcal{H}$

Then

$$\{f : f \text{ nonnegative } \sigma(\mathcal{A})\text{-measurable}\} \subset \mathcal{H}$$

## 2. 集函数与测度

**2.1. 集函数.**  $\mathcal{E}$  is a collection of subsets of  $E$ .

**Def 2.** *Set function,  $\mu : \mathcal{E} \mapsto \mathbb{R} \cup \{\pm\infty\}$ .*

**Def 3.** *Nonnegative set function,  $\mu : \mathcal{E} \mapsto \mathbb{R} \cup \{\infty\}$ .*

**Def 4.**  *$\mu$  is finite if,  $\forall A \in \mathcal{E}, |\mu(A)| < \infty$ .*

**Def 5.**  *$\mu$  is  $\sigma$ -finite on  $\mathcal{E}$  if,  $\forall A \in \mathcal{E}$ , there exist  $\{A_n\} \subset \mathcal{E}$ ,  $A = \bigcup_n A_n$  with  $|\mu(A_n)| < \infty$ .*

**Def 6.**  *$\mu$  is additive if,  $\forall A, B \in \mathcal{E}, AB = \emptyset$ ,*

$$\mu(A + B) = \mu(A) + \mu(B).$$

**Def 7.**  *$\mu$  is countably additive if,  $\forall A_i \in \mathcal{E}, i = 1, 2, \dots$ , disjoint,*

$$\mu\left(\sum_i A_i\right) = \sum_i \mu(A_i).$$

**Def 8.**  $\emptyset \in \mathcal{E}$ .  $\mu$  is a measure on  $\mathcal{E}$  if it is nonnegative, countably additive,  $\mu(\emptyset) = 0$ .

**Example 1.**  $(X, \mathcal{F})$  measurable space,  $x \in X$ ,

$$\delta_x(A) = 1_A(x), \quad \forall A \in \mathcal{F}.$$

$$x_1, \dots, x_n \in X,$$

$$\mu(A) = \sum_i \delta_{x_i}(A), \quad \forall A \in \mathcal{F}.$$

**Example 2.**  $F$  real-valued nonnegative, non-decreasing, right continuous. Semi-ring on  $\mathbb{R}$ ,

$$\mathcal{A} = \{(a, b] : a, b, \in \mathbb{R}\}.$$

Then

$$\mu((a, b]) = F(b) - F(a)$$

defines a measure  $\mu$ . It is unique on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ .

PROOF. **1.** Additivity.  $(a_i, b_i]$ ,  $i = 1, \dots, n$ , disjoint,  $(a, b] = \bigcup_i^n (a_i, b_i]$ , then

$$\mu((a, b]) = \sum_{i=1}^n \mu((a_i, b_i]).$$

**2.**  $(a_i, b_i]$ ,  $i = 1, \dots$ , disjoint,  $\bigcup_i (a_i, b_i] \subset (a, b]$ , then

$$\sum_{i=1}^{\infty} \mu((a_i, b_i]) \leq \mu((a, b]).$$

**3.**  $(a_i, b_i]$ ,  $i = 1, \dots, n$ ,  $(a, b] \subset \bigcup_i^n (a_i, b_i]$ , then

$$\mu((a, b]) \leq \sum_{i=1}^n \mu((a_i, b_i]).$$

4.  $(a_i, b_i]$ ,  $i = 1, \dots$ , disjoint,  $\bigcup_i (a_i, b_i] = (a, b]$ , then

$$\mu((a, b]) = \sum_{i=1}^{\infty} \mu((a_i, b_i]).$$

$\forall \varepsilon > 0$ , there is  $\delta_i > 0$ ,

$$F(b_i + \delta_i) - F(b_i) < \frac{\varepsilon}{2^i}.$$

$\forall \theta > 0$ ,  $\{(a_i, b_i + \delta_i) : i\}$  is an open cover of  $[a + \theta, b]$ , there exists  $n_0$

$$(a + \theta, b] \subset \bigcup_i^{n_0} (a_i, b_i + \delta_i].$$



By **3.**,

$$\begin{aligned}\mu((a + \theta, b]) &\leq \sum_{i=1}^{n_0} \mu((a_i, b_i + \delta_i]) \\ &= \sum_{i=1}^{n_0} (F(b_i + \delta_i) - F(b_i)) \\ &\leq \sum_{i=1}^{n_0} (F(b_i) - F(b_i)) + \sum_{i=1}^{n_0} \frac{\varepsilon}{2^i} \\ &\leq \sum_{i=1}^{\infty} (F(b_i) - F(b_i)) + \varepsilon.\end{aligned}$$

□

**2.2. 半环上非负集函数.**  $\mathcal{E}$  is a collection of subsets of  $E$ ,  $\mu$  is a nonnegative set function on  $\mathcal{E}$ .

**Def 9.** *Monotonicity:*  $\forall A \subset B \in \mathcal{E}$ ,

$$\mu(A) \leq \mu(B).$$

**Def 10.** *Countably subadditive:*  $\forall A_i \in \mathcal{E}, i = 1, 2, \dots, \bigcup_{i=1}^{\infty} A_i \in \mathcal{E}$ ,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mu(A_i).$$

**Def 11.** *Continuity from below:*  $A_i \in \mathcal{E}, A_i \uparrow A \in \mathcal{E}$ ,

$$\lim_n \mu(A_i) = \mu(A).$$

**Def 12.** *Continuity from above:*  $A_i \in \mathcal{E}, A_i \downarrow A \in \mathcal{E}, \mu(A_1) < \infty$ ,

$$\lim_n \mu(A_i) = \mu(A).$$

REMARK 1. **Note** finiteness is part of the definition of continuity from above.

$\mathcal{S}$  is a semi-ring on  $E$ ,  $\mu$  is a nonnegative set function on  $\mathcal{S}$ .

Suppose  $\mu$  is **additive**.

1.  $\mu(\emptyset) = 0, +\infty$ .

PROOF.  $\emptyset \in \mathcal{S}$ . By additivity

$$\mu(\emptyset) = \sum_{i=1}^n \mu(\emptyset).$$

$\mu(\emptyset)$  equals 0, or  $\infty$ .

□

2. Monotonicity.

PROOF.  $A, B \in \mathcal{S}$ ,  $A \subset B$ . There exist disjoint  $C_1, \dots, C_k \in \mathcal{S}$ ,

$$B \setminus A = \bigcup_{i=1}^k C_i.$$

$$B = A \cup (B \setminus A) = A \cup \left( \bigcup_{i=1}^k C_i \right).$$

By additivity

$$\mu(B) = \mu(A) + \sum_{i=1}^k \mu(C_i) \geq \mu(A).$$

□

Suppose  $\mu$  is **countably additive**.

**3.** Continuity from below.

PROOF.  $A_i \in \mathcal{S}$ ,  $A_i \uparrow A \in \mathcal{S}$ . There exist disjoint  $C_{n,1}, \dots, C_{n,k_n} \in \mathcal{S}$ ,

$$B_n \triangleq A_n \setminus A_{n-1} = \bigcup_{i=1}^{k_n} C_{n,i}.$$

$$(A_0 = \emptyset)$$

$$\begin{aligned}\mu(A) &= \mu\left(\bigcup_{n=1}^{\infty} A_n\right) = \mu\left(\bigcup_{n=1}^{\infty} B_n\right) = \mu\left(\bigcup_{n=1}^{\infty} \bigcup_{i=1}^{k_n} C_{n,i}\right) \\ &= \sum_{n=1}^{\infty} \sum_{i=1}^{k_n} \mu(C_{n,i}) = \lim_N \sum_{n=1}^N \sum_{i=1}^{k_n} \mu(C_{n,i}) \\ &= \lim_N \mu\left(\bigcup_{n=1}^N \bigcup_{i=1}^{k_n} C_{n,i}\right) = \lim_n \mu(A_n).\end{aligned}$$

□

4. Continuity from above.

PROOF. (**WRONG PROOF**)  $A_i \in \mathcal{S}$ ,  $A_i \downarrow A \in \mathcal{S}$ ,  $\mu(A_1) < \infty$ .  
Clearly

$$\mu\left(\bigcap_{n=1}^{\infty} A_n\right) \leq \mu(A_i) \leq \mu(A_1) < \infty.$$

$$\lim_n \mu(A_n) = \mu\left(\bigcap_{n=1}^{\infty} A_n\right)$$

$$\iff$$

$$\mu(A_1) - \lim_n \mu(A_n) = \mu(A_1) - \mu\left(\bigcap_{n=1}^{\infty} A_n\right)$$

$$\iff$$

$$\lim_n \mu(A_1 \setminus A_n) = \mu\left(A_1 \setminus \bigcap_{n=1}^{\infty} A_n\right) = \mu\left(\bigcup_{n=1}^{\infty} (A_1 \setminus A_n)\right).$$

□

## 5. Subadditivity.

PROOF. Analogous to continuity from below. □

### 2.3. 环上非负集函数.

**Thm 6.**  $\mathcal{R}$  is a ring.  $\mu$  is nonnegative additive.

(1)  $\mu$  countably additive



(2)  $\mu$  countably subadditive



(3)  $\mu$  continuity from below



(4)  $\mu$  continuity from above



(5)  $\mu$  continuity from above at  $\emptyset$ .

If  $\mu$  is finite, (5) implies (1).

PROOF. **1.** Already have:  $(1) \implies (2)$ ,  $(1) \implies (3)$ ,  $(1) \implies (4)$ ,  $(4) \implies (5)$ .

**2.**  $(2) \implies (1)$ . Suppose  $A_i \in \mathcal{R}$ ,  $i = 1, 2, \dots$ , disjoint,  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{R}$ .

By countable subadditivity,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mu(A_i).$$

By monotonicity and additivity,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) \geq \mu\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mu(A_i), \quad \forall n.$$

Sending  $n \rightarrow \infty$ ,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) \geq \sum_{i=1}^{\infty} \mu(A_i).$$



3. (3)  $\implies$  (1). Suppose  $A_i \in \mathcal{R}$ ,  $i = 1, 2, \dots$ , disjoint,  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{R}$ .

Since

$$\bigcup_{i=1}^n A_i \uparrow \bigcup_{i=1}^{\infty} A_i,$$

by continuity from below,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_n \mu\left(\bigcup_{i=1}^n A_i\right) = \lim_n \sum_{i=1}^n \mu(A_i) = \sum_{i=1}^{\infty} \mu(A_i).$$

4. (5)  $\implies$  (1). Suppose  $A_i \in \mathcal{R}$ ,  $i = 1, 2, \dots$ , disjoint,  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{R}$ .

Then,  $\forall n$ ,

$$\bigcup_{i=1}^n A_i \in \mathcal{R} \text{ and } \bigcup_{i=n+1}^{\infty} A_i = \bigcup_{i=1}^{\infty} A_i \setminus \bigcup_{i=1}^n A_i \in \mathcal{R}.$$

By additivity

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \mu\left(\bigcup_{i=1}^n A_i\right) + \mu\left(\bigcup_{i=n+1}^{\infty} A_i\right).$$

Since  $\mu$  is finite

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) < \infty.$$

The continuity from above at  $\emptyset$  yields,

$$\lim_n \mu\left(\bigcup_{i=n+1}^{\infty} A_i\right) = 0.$$

Hence

$$\begin{aligned}\mu\left(\bigcup_{i=1}^{\infty} A_i\right) &= \lim_n \mu\left(\bigcup_{i=1}^n A_i\right) + \lim_n \mu\left(\bigcup_{i=n+1}^{\infty} A_i\right) \\ &= \lim_n \sum_{i=1}^n \mu(A_i) = \sum_{i=1}^{\infty} \mu(A_i).\end{aligned}$$

□

### 3. Carathéodory's 延拓

#### 3.1. 外测度.

**Def 13.**  $\mu^*$  is an outer measure on  $E$  if

(1)  $\mu^*(\emptyset) = 0$

(2)  $\forall A, B \in 2^E$ , if  $A \subset B$ , then

$$\mu^*(A) \leq \mu^*(B)$$

(3) If  $A_i \in 2^E, i = 1, 2, \dots$ ,

$$\mu^*\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mu^*(A_i)$$

**Thm 7.** Let  $\mathcal{E}$  be a collection of sets on  $E$ ,  $\emptyset \in \mathcal{E}$ .  $\mu$  is a nonnegative set function on  $\mathcal{E}$  with  $\mu(\emptyset) = 0$ . Define,  $\forall A \in 2^E$ ,

$$\mu^*(A) = \inf \left\{ \sum_{i=1}^{\infty} \mu(A_i) : A_i \in \mathcal{E}, A \subset \bigcup_{i=1}^{\infty} A_i \right\}.$$

Then  $\mu^*(A)$  is an outer measure.

PROOF. **1.**  $\mu^*(\emptyset) = 0$  since  $\emptyset \in \mathcal{E}$ ,  $\emptyset \subset \bigcup_{i=1}^{\infty} \emptyset$ .

**2.** If  $A \subset B$ ,  $B \subset \bigcup_{i=1}^{\infty} B_i$ , then  $A \subset \bigcup_{i=1}^{\infty} B_i$ , from the definition  $\mu^*(A) \leq \mu^*(B)$ .

**3.** Let  $A_i \in 2^E, i = 1, 2, \dots, \varepsilon > 0$ . There are  $A_{i,k} \in \mathcal{E}$ ,  $A_i \subset \bigcup_{k=1}^{\infty} A_{i,k}$ ,

$$\sum_{k=1}^{\infty} \mu(A_{i,k}) \leq \mu^*(A_i) + \frac{\varepsilon}{2^i}, \quad \forall i.$$

Since

$$\bigcup_{i=1}^{\infty} A_i \subset \bigcup_{i=1}^{\infty} \bigcup_{k=1}^{\infty} A_{i,k},$$

$$\begin{aligned}
\mu^*\left(\bigcup_{i=1}^{\infty} A_i\right) &\leq \sum_{i=1}^{\infty} \sum_{k=1}^{\infty} \mu(A_{i,k}) \\
&\leq \sum_{i=1}^{\infty} \left[ \mu^*(A_i) + \frac{\varepsilon}{2^i} \right] \leq \sum_{i=1}^{\infty} \mu^*(A_i) + \varepsilon.
\end{aligned}$$

□

**Def 14.**  $\mu^*$  is an outer measure on  $E$ .  $A \in 2^E$  is  $\mu^*$ -measurable if  

$$\mu^*(D) = \mu^*(D \cap A) + \mu^*(D \cap A^c), \quad \forall D \in 2^E.$$

The class of  $\mu^*$ -measurable sets is denoted by  $\mathcal{F}_\mu^*$ .

**Def 15.** Let  $\mu$  be a measure on a  $\sigma$ -field  $\mathcal{F}$  of  $E$ , the measure space  $(E, \mathcal{F}, \mu)$  is complete if

$$A \in \mathcal{F}, \quad \mu(A) = 0 \implies B \in \mathcal{F}, \quad \forall B \subset A.$$

**Thm 8** (Carathéodory). *Let  $\mathcal{E}$  be a collection of sets on  $E$ ,  $\emptyset \in \mathcal{E}$ .  $\mu$  is a nonnegative set function on  $\mathcal{E}$  with  $\mu(\emptyset) = 0$ .*

(1)  $\mathcal{F}_\mu^*$  is a  $\sigma$ -field.

(2)  $(E, \mathcal{F}_\mu^*, \mu^*)$  is a complete measure space.

PROOF. 1. Obviously,  $E \in \mathcal{F}_\mu^*$  and  $A^c \in \mathcal{F}_\mu^*$  if  $A \in \mathcal{F}_\mu^*$ .

2. If  $A_1, A_2 \in \mathcal{F}_\mu^*$ , then  $A_1 \cup A_2, A_1 \cap A_2 \in \mathcal{F}_\mu^*$ .

$\forall D \in 2^E$ , we note

$$D \cap (A_1 \cup A_2) = (D \cap A_1) \cup (D \cap A_1^c \cap A_2).$$

Then

$$\begin{aligned} & \mu^*(D \cap (A_1 \cup A_2)) + \mu^*(D \cap (A_1 \cup A_2)^c) \\ & \leq \mu^*(D \cap A_1) + \mu^*(D \cap A_1^c \cap A_2) + \mu^*(D \cap A_1^c \cap A_2^c) \quad (\text{subadditivity}) \\ & \leq \mu^*(D \cap A_1) + \mu^*(D \cap A_1^c) \quad (A_2 \in \mathcal{F}_\mu^*) \\ & = \mu^*(D) \quad (A_1 \in \mathcal{F}_\mu^*). \end{aligned}$$

Hence

$$A_1 \cup A_2 \in \mathcal{F}_\mu^*.$$

It follows that

$$(A_1 \cap A_2)^c = A_1^c \cup A_2^c \in \mathcal{F}_\mu^*.$$

**3. Finite additivity.** If  $A_1, \dots, A_n \in \mathcal{F}_\mu^*$  disjoint, then  $\forall D \in 2^E$ ,

$$\mu^* \left( D \cap \left( \bigcup_{i=1}^n A_i \right) \right) = \sum_{i=1}^n \mu^*(D \cap A_i).$$



Indeed, since  $A_1 \in \mathcal{F}_\mu^*$ ,

$$\begin{aligned}
& \mu^* \left( D \cap \left( \bigcup_{i=1}^n A_i \right) \right) \\
&= \mu^* \left( D \cap \left( \bigcup_{i=1}^n A_i \right) \cap A_1 \right) + \mu^* \left( D \cap \left( \bigcup_{i=1}^n A_i \right) \cap A_1^c \right) \\
&= \mu^*(D \cap A_1) + \mu^* \left( D \cap \left( \bigcup_{i=2}^n A_i \right) \right) = \cdots = \sum_{i=1}^n \mu^*(D \cap A_i)
\end{aligned}$$

4. If  $A_1, A_2, \dots \in \mathcal{F}_\mu^*$ , then  $A \triangleq \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}_\mu^*$ .

We can assume that  $A_1, A_2, \dots \in \mathcal{F}_\mu^*$  are disjoint. Indeed, by **1** and **2**,  $B_i = A_i \setminus \left( \bigcup_{j < i} A_j \right) \in \mathcal{F}_\mu^*$ , are disjoint and  $\bigcup_{i=1}^n B_i = \bigcup_{i=1}^n A_i$ ,

$\forall n$ . Let

$$C_n = \bigcup_{i=1}^n A_i \in \mathcal{F}_\mu^*, \quad \forall n.$$

Since  $A_1, A_2, \dots$  are disjoint, we can use **3** (the finite additivity).  $\forall D \in 2^E$ ,

$$\begin{aligned} \mu^*(D) &= \mu^*(D \cap C_n) + \mu^*(D \cap C_n^c) \\ &= \sum_{i=1}^n \mu^*(D \cap C_i) + \mu^*(D \cap C_n^c) \\ &\geq \sum_{i=1}^n \mu^*(D \cap C_i) + \mu^*(D \cap A^c), \quad \forall n. \end{aligned}$$

Let  $n \rightarrow \infty$ , note  $A \subset \bigcup_{i=1}^{\infty} C_i$  and use subadditivity of outer measure

$$\mu^*(D) \geq \sum_{i=1}^{\infty} \mu^*(D \cap C_i) + \mu^*(D \cap A^c) \geq \mu^*(D \cap A) + \mu^*(D \cap A^c).$$

## 5. Countable additivity.

If  $A_1, A_2, \dots, \in \mathcal{F}_\mu^*$  are disjoint, use **3** and send  $n \rightarrow \infty$ ,

$$\mu^*\left(\bigcup_{i=1}^{\infty} A_i\right) \geq \mu^*\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mu^*(A_i), \quad \forall n.$$

The opposite inequality is subadditivity of outer measure.

**6. Completeness.** If  $A \in \mathcal{F}_\mu^*$ ,  $\mu^*(A) = 0$  and  $B \subset A$ , then  $\mu^*(B) = 0$ .  $\forall D \in 2^E$ ,

$$\mu^*(D) \geq \mu^*(D \cap B^c) = \mu^*(D \cap B) + \mu^*(D \cap B^c).$$

So  $B \in \mathcal{F}_\mu^*$ . □

## 3.2. 域上测度的延拓.

**Thm 9.** *If  $\mu$  is a measure on a field  $\mathcal{A}$  with the generated outer measure  $\mu^*$ . Then*

(1)  $\mathcal{A} \subset \mathcal{F}_\mu^*$  thus  $\sigma(\mathcal{A}) \subset \mathcal{F}_\mu^*$ .

(2)  $\mu^*$  is an extension of  $\mu$  to  $\sigma(\mathcal{A})$  in the sense that

$$\mu(A) = \mu^*(A), \quad \forall A \in \mathcal{A}.$$

PROOF. 1. Let  $A \subset \mathcal{A}$ . If  $A_i \in \mathcal{A}$ ,  $A \subset \bigcup_{i=1}^{\infty} A_i$ , then

$$(3.1) \quad \mu(A) \leq \sum_{i=1}^{\infty} \mu(A_i).$$

Indeed,

$$\mu\left(A \cap \bigcup_{i=1}^n A_i\right) \leq \mu\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n \mu(A_i).$$

Let  $n \rightarrow \infty$  and use that  $\mu$  is a measure to get (3.1). So

$$\mu(A) \leq \mu^*(A).$$

Since  $A \subset \mathcal{A}$ ,  $A_1 = A$ ,  $A_2 = A_3 \dots = \emptyset$  form a countable cover of  $A$ , so

$$\mu^*(A) \leq \mu(A).$$

**2.** Fix  $A \subset \mathcal{A}$ , will prove  $A \in \mathcal{F}_\mu^*$ .  $\forall D \in 2^E$ , it is enough to show that

$$\mu^*(D) \geq \mu^*(D \cap A) + \mu^*(D \cap A^c).$$

There is nothing to prove if  $\mu^*(D) = \infty$ , so we assume that  $\mu^*(D) < \infty$ . Then,  $\forall \varepsilon > 0$ , there exist  $A_i \in \mathcal{A}$ ,  $D \subset \bigcup_{i=1}^{\infty} A_i$  so that

$$\sum_{i=1}^{\infty} \mu(A_i) \leq \mu^*(D) + \varepsilon.$$

Since  $\mathcal{A}$  is a field,

$$A_i \cap A, A_i \cap A^c \in \mathcal{A}.$$

By **1** and the additivity of  $\mu$ ,

$$\begin{aligned}\mu(A_i) &= \mu(A_i \cap A) + \mu(A_i \cap A^c) \\ &= \mu^*(A_i \cap A) + \mu^*(A_i \cap A^c).\end{aligned}$$

Summing over  $i$  gives

$$\begin{aligned}\sum_{i=1}^{\infty} \mu(A_i) &= \sum_{i=1}^{\infty} \mu^*(A_i \cap A) + \sum_{i=1}^{\infty} \mu^*(A_i \cap A^c) \\ &\geq \mu^*(D \cap A) + \mu^*(D \cap A^c).\end{aligned}$$

So

$$\mu^*(D) + \varepsilon \geq \sum_{i=1}^{\infty} \mu(A_i) \geq \mu^*(D \cap A) + \mu^*(D \cap A^c).$$

□

**Thm 10** (Uniqueness). *Let  $\mathcal{P}$  be a  $\pi$ -system on  $E$ ,  $\mu$  and  $\nu$  measures on  $\sigma(\mathcal{P})$ . Assume that*

*(1)  $\mu$  and  $\nu$  agree on  $\mathcal{P}$ .*

(2) There are  $B_i \in \mathcal{P}$ ,  $i = 1, 2, \dots$ , disjoint so that  $\bigcup_{i=1}^{\infty} B_i = E$  and

$$\mu(B_i) < \infty.$$

Then  $\mu$  and  $\nu$  are equal on  $\sigma(\mathcal{P})$ .

PROOF. 1. Let  $B \in \mathcal{P}$  have  $\mu(B) < \infty$ . Define

$$\mathcal{L} = \{A \in \sigma(\mathcal{P}) : \mu(A \cap B) = \nu(A \cap B)\}.$$

$\mathcal{L}$  is a  $\lambda$ -system (finiteness is needed to justify sets subtraction!),  $\mathcal{P} \subset \mathcal{L}$ . So

$$\sigma(\mathcal{P}) \subset \mathcal{L},$$

i.e.

$$\mu(A \cap B) = \nu(A \cap B), \quad \forall A \in \sigma(\mathcal{P}).$$

2.  $\forall A \in \sigma(\mathcal{P})$ , use (2) to write it as disjoint union,

$$A = \bigcup_{i=1}^{\infty} (A \cap B_i), \quad \mu(A \cap B_i) \leq \mu(B_i) < \infty.$$

Then, by 1,

$$\begin{aligned}\mu(A) &= \mu\left(\bigcup_{i=1}^{\infty}(A \cap B_i)\right) = \sum_{i=1}^{\infty} \mu(A \cap B_i) \\ &= \sum_{i=1}^{\infty} \nu(A \cap B_i) = \nu\left(\bigcup_{i=1}^{\infty}(A \cap B_i)\right) = \nu(A).\end{aligned}$$

□

▷ 2. The condition Theorem 10 (2) can be replaced with either one of the following:

(2')  $\mathcal{P}$  is a semi-ring,  $E \in \mathcal{P}$  and  $\mu$  is  $\sigma$ -finite on  $\mathcal{P}$ .

(2'') there are  $B_1, B_2, \dots \in \mathcal{P}$ , so that  $B_i \uparrow E$  and  $\mu(B_i) < \infty$ .

### 3.3. 半环上测度的延拓.

**Thm 11.** Let  $\mu$  be a measure on the semi-ring  $\mathcal{S}$  with the generated outer measure  $\mu^*$ . Then

(1)  $\mathcal{S} \subset \mathcal{F}_{\mu}^*$  thus  $\sigma(\mathcal{S}) \subset \mathcal{F}_{\mu}^*$ .



(2)  $\mu^*$  is an extension of  $\mu$  to  $\sigma(\mathcal{S})$  in the sense that

$$(3.2) \quad \mu(A) = \mu^*(A), \quad \forall A \in \mathcal{S}.$$

(3) Assume that there are  $B_i \in \mathcal{S}$ ,  $i = 1, 2, \dots$ , disjoint so that  $\bigcup_{i=1}^n B_i = E$  and  $\mu(B_i) < \infty$ , then the extension of  $\mu$  to  $\sigma(\mathcal{S})$  is unique.

PROOF. Let  $\bar{\mu}$  be the outer measure generated by  $\mu$ .

1.  $\bar{\mu}$  agrees with  $\mu$  on  $\mathcal{S}$ .

The proof is identical to Theorem 9 (1).

2. Fix  $A \subset \mathcal{S}$ , will prove  $A \in \mathcal{F}_\mu^*$ .

The proof is identical to Theorem 9 (2). The difference is  $A_i \cap A^c$  is replaced with disjoint union of sets in  $\mathcal{S}$ .

3. Uniqueness. Apply Theorem 10 to conclude. □

### 3.4. Approximating $\mu^*|_{\mathcal{F}_\mu^*}$ by $\mu^*|_{\sigma(\mathcal{S})}$ .

**Thm 12.** *Let  $\mu$  be a measure on the semi-ring  $\mathcal{S}$  with the generated outer measure  $\mu^*$ . Suppose  $E \in \mathcal{S}$ .*

(1)  *$\forall A \in \mathcal{F}_\mu^*$ , there is  $B \in \sigma(\mathcal{S})$  such that  $A \subset B$  and*

$$\mu^*(A) = \mu^*(B).$$

(2) *If  $\mu$  is  $\sigma$ -finite on  $\mathcal{S}$ , then  $\forall A \in \mathcal{F}_\mu^*$ , there is  $B \in \sigma(\mathcal{S})$  such that  $A \subset B$  and*

$$\mu^*(B \setminus A) = 0.$$

PROOF.

1. There is nothing to prove if  $\mu^*(A) = \infty$ , we assume that  $\mu^*(A) < \infty$ . There are  $B_{n,i} \in \mathcal{S}$ ,  $A \subset \bigcup_{i=1}^{\infty} B_{n,i}$ ,

$$\sum_{i=1}^{\infty} \mu(B_{n,i}) < \mu^*(A) + \frac{1}{n}.$$

Set

$$B = \bigcap_{n=1}^{\infty} \bigcup_{i=1}^{\infty} B_{n,i}.$$

Then  $A \subset B \in \sigma(\mathcal{S})$ ,

$$\mu^*(A) \leq \mu^*(B).$$

Moreover

$$\mu^*(B) \leq \mu^*\left(\bigcup_{i=1}^{\infty} B_{n,i}\right) \leq \sum_{i=1}^{\infty} \mu(B_{n,i}) \leq \mu^*(A) + \frac{1}{n}.$$

It follows that

$$\mu^*(B) \leq \mu^*(A).$$

**2.** If  $\mu$  is *finite* on  $\mathcal{S}$ , then by **1**,  $\forall A \in \mathcal{F}_{\mu}^*$ , there is  $B \in \sigma(\mathcal{S})$  such that  $A \subset B$  and

$$\mu^*(A) = \mu^*(B).$$

Since  $\mu^*$  is a measure on  $\mathcal{F}_{\mu}^*$ , this gives

$$\mu^*(B \setminus A) = 0.$$

The  $\sigma$ -finite case follows from similar argument as in step **3** of Theorem 11.  $\square$

### 3.5. Approximating $\mu|_{\sigma(\mathcal{A})}$ by $\mu|_{\mathcal{A}}$ .

**Thm 13.** *Let  $\mu$  be a measure on the field  $\mathcal{A}$  with the generated outer measure  $\mu^*$ . For any  $A \in \sigma(\mathcal{A})$  with  $\mu^*(A) < \infty$ ,  $\forall \varepsilon > 0$ , there is  $B \in \mathcal{A}$  such that  $\mu^*(A \Delta B) < \varepsilon$ .*

If, in the last Theorem, the measure  $\mu$  is defined on  $\sigma(\mathcal{A})$  and  $\sigma$ -finite on  $\mathcal{A}$ , then  $\mu$  must equal  $\mu^*$  on  $\sigma(\mathcal{A})$  by uniqueness, we can use  $\mu$  in place of  $\mu^*$  in the conclusion.

**Thm 14.** *Let  $\mathcal{A}$  be a field,  $\mu$  a measure on  $\sigma(\mathcal{A})$  and  $\sigma$ -finite on  $\mathcal{A}$ . For any  $A \in \sigma(\mathcal{A})$  with  $\mu(A) < \infty$ ,  $\forall \varepsilon > 0$ , there is  $B \in \mathcal{A}$  such that  $\mu(A \Delta B) < \varepsilon$ .*

### 3.6. Completion of a measure space.

**Thm 15.** *Let  $(X, \mathcal{F}, \mu)$  be a measure space,*

$$\bar{\mathcal{F}} \triangleq \{A \cup N : A \in \mathcal{F}, N \subset B \text{ for some } B \in \mathcal{F} \text{ with } \mu(B) = 0\}.$$

Define

$$\bar{\mu}(A \cup N) = \mu(A), \quad \forall A \in \bar{\mathcal{F}}.$$

Then  $(X, \bar{\mathcal{F}}, \bar{\mu})$  is a complete measure space.

Clearly the Theorem says

$$\bar{\mu}(A) = \mu(A), \quad \forall A \in \bar{\mathcal{F}}.$$

PROOF. 1.  $\bar{\mathcal{F}}$  is a  $\sigma$ -field.

Suppose  $A \cup N \in \bar{\mathcal{F}}$  where  $A \in \mathcal{F}$ ,  $N \subset B$ ,  $B \in \mathcal{F}$  with  $\mu(B) = 0$ .  
Then

$$(A \cup N)^c = (A^c \cap B^c) \cup (B \cap A^c \cap N^c) \in \bar{\mathcal{F}}.$$

Suppose  $A_i \cup N_i \in \bar{\mathcal{F}}$  where  $A_i \in \mathcal{F}$ ,  $N_i \subset B_i$ ,  $B_i \in \mathcal{F}$  with  $\mu(B_i) = 0$ . Then

$$\bigcup_{i=1}^{\infty} (A_i \cup N_i) = \left( \bigcup_{i=1}^{\infty} A_i \right) \cup \left( \bigcup_{i=1}^{\infty} N_i \right) \in \bar{\mathcal{F}},$$

since

$$\bigcup_{i=1}^{\infty} N_i \subset \bigcup_{i=1}^{\infty} B_i \in \mathcal{F}$$

and

$$\mu\left(\bigcup_{i=1}^{\infty} B_i\right) = \sum_{i=1}^{\infty} \mu(B_i) = 0.$$

**2.** The definition of  $\bar{\mu}$  nonambiguous, i.e.

$$A_1 \cup N_1 = A_2 \cup N_2 \in \tilde{\mathcal{F}} \implies \bar{\mu}(A_1 \cup N_1) = \bar{\mu}(A_2 \cup N_2).$$

Here  $N_i \subset B_i$  for some  $B_i \in \mathcal{F}$  with  $\mu(B_i) = 0$ ,  $i = 1, 2$ .

$$\bar{\mu}(A_1 \cup N_1) = \mu(A_1) = \mu(A_1 \cup B_1 \cup B_2) \geq \mu(A_2) = \bar{\mu}(A_2 \cup N_2).$$

By symmetry,

$$\bar{\mu}(A_1 \cup N_1) \leq \bar{\mu}(A_2 \cup N_2).$$

(In fact

$$A_1 \cup B_1 \cup B_2 = A_1 \cup N_1 \cup B_1 \cup B_2 = A_2 \cup N_2 \cup B_1 \cup B_2 = A_2 \cup B_1 \cup B_2$$

so

$$\mu(A_1 \cup B_1 \cup B_2) = \mu(A_2).$$

)

**3. Countable additivity.** Suppose  $A_i \cup N_i \in \bar{\mathcal{F}}$  disjoint, where  $A_i \in \mathcal{F}$ ,  $N_i \subset B_i$ ,  $B_i \in \mathcal{F}$  with  $\mu(B_i) = 0$ . Then

$$\bar{\mu}\left(\bigcup_{i=1}^{\infty} (A_i \cup N_i)\right) = \mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i) = \sum_{i=1}^{\infty} \bar{\mu}(A_i \cup N_i).$$

**4. Completeness.** Let  $A \cup N \in \bar{\mathcal{F}}$ ,  $N \subset B$ ,  $B \in \mathcal{F}$  with  $\mu(B) = 0$  and  $\bar{\mu}(A \cup N)$ , then

$$\mu(A \cup B) = \mu(A) = \bar{\mu}(A \cup N) = 0.$$

So for any  $C \subset A \cup N$ ,  $C \subset A \cup B$ ,

$$C = \emptyset \cup C \in \bar{\mathcal{F}}.$$

□

**Thm 16.** *Suppose that  $\mu$  is  $\sigma$ -finite on the semi-ring  $\mathcal{S}$  with the generated outer measure  $\mu^*$ . Then  $(X, \mathcal{F}_\mu^*, \mu^*)$  is the completion of  $(X, \sigma(\mathcal{S}), \mu^*)$ .*

PROOF. Let

$$\bar{\mathcal{F}} \triangleq \{A \cup N : A \in \sigma(\mathcal{S}), N \subset B \text{ for some } B \in \sigma(\mathcal{S}) \text{ with } \mu(B) = 0\}.$$

It is enough to show that

$$\mathcal{F}_\mu^* = \bar{\mathcal{F}}.$$

Since  $(X, \mathcal{F}_\mu^*, \mu^*)$  is a complete measure space,

$$\bar{\mathcal{F}} \subset \mathcal{F}_\mu^*.$$

Let  $A \in \mathcal{F}_\mu^*$ , by Theorem 12 there exist  $B, C \in \sigma(\mathcal{S})$  so that

$$A \subset B, \mu^*(B \setminus A) = 0; B \setminus A \subset C, \mu^*(C) = \mu^*(B \setminus A) = 0.$$

Writing

$$A = (B \cap C^c) \cup (A \cap C),$$

we get that  $B \cap C^c \in \sigma(\mathcal{S})$ ,  $(A \cap C) \subset C$ ,  $\mu^*(C) = 0$ , so  $A \in \bar{\mathcal{F}}$ .  $\square$



## 4. 收敛

**4.1. 可测函数的收敛.**  $(E, \mathcal{F}, \mu)$  a measure space,  $f_n \in \mathcal{F}$ ,  $i = 1, 2, \dots$ ,  $f \in \mathcal{F}$

**Def 16.** *Almost everywhere convergence,  $f_n \xrightarrow{a.e.} f$ :*

$$\mu\left(\lim_n f_n \neq f\right) = 0.$$

**Def 17.** *Convergence in measure,  $f_n \xrightarrow{\mu} f$ :  $\forall \varepsilon > 0$ ,*

$$\lim_n \mu(|f_n - f| > \varepsilon) = 0.$$

Evidently

$$\begin{aligned} f_n \xrightarrow{a.e.} f &\iff \forall \varepsilon > 0, \mu\left(\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} \{|f_m - f| > \varepsilon\}\right) = 0 \\ &\iff \forall \varepsilon > 0, \mu(\{|f_n - f| > \varepsilon\} \text{ i.o.}) = 0. \end{aligned}$$

Recall

$$x \in \limsup A_n \iff x \in A_n \text{ i.o.}$$

**Thm 17.** *If  $\mu$  is finite, then*

$$f_n \xrightarrow{a.e.} f \implies f_n \xrightarrow{\mu} f.$$

PROOF. Indeed,

$$\mu(|f_n - f| > \varepsilon) \leq \mu\left(\bigcup_{m=n}^{\infty} \{|f_m - f| > \varepsilon\}\right), \quad \forall n.$$

Let  $n \rightarrow \infty$  and use continuity from above (requires finiteness of  $\mu$ )

$$\begin{aligned} \limsup_n \mu(|f_n - f| > \varepsilon) &\leq \lim_n \mu\left(\bigcup_{m=n}^{\infty} \{|f_m - f| > \varepsilon\}\right) \\ &= \mu\left(\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} \{|f_m - f| > \varepsilon\}\right) = 0. \end{aligned}$$

(or use

$$\limsup_n \mu(A_n) \leq \mu\left(\limsup_n A_n\right).$$

)

□

**Def 18.** *Almost uniform convergence,  $f_n \xrightarrow{a.u.} f$ :  $\forall \varepsilon > 0$ , there is  $A_\varepsilon \in \mathcal{F}$  so that  $\mu(A_\varepsilon) < \varepsilon$ ,*

$$\lim_n \sup_{x \notin A_\varepsilon} |f_n - f| = 0.$$

Compare with Egoroff's Theorem on *finite* measure!

**Thm 18.**  $f_n \xrightarrow{a.u.} f$  if and only if  $\forall \varepsilon > 0$ ,

$$\lim_n \mu \left( \bigcup_{m=n}^{\infty} \{|f_m - f| > \varepsilon\} \right) = 0.$$

PROOF. 1. " $\implies$ ".  $\forall \varepsilon > 0$ , there is  $A_\varepsilon$  so that  $\mu(A_\varepsilon) < \varepsilon$  and

$$\lim_m \sup_{x \notin A_\varepsilon} |f_m - f| = 0.$$

So,  $\forall \varepsilon' > 0$ , there is  $n_0 \in \mathbb{N}$  such that

$$\sup_{x \notin A_\varepsilon} |f_m - f| \leq \varepsilon', \quad \forall m \geq n_0.$$

This translates to

$$\bigcup_{m=n_0}^{\infty} \{|f_m - f| > \varepsilon'\} \subset A_{\varepsilon}.$$

Therefore

$$\mu\left(\bigcup_{m=n_0}^{\infty} \{|f_m - f| > \varepsilon'\}\right) \leq \mu(A_{\varepsilon}) < \varepsilon.$$

**2.** "  $\Leftarrow$  ".  $\forall \varepsilon > 0$  and  $k \in \mathbb{N}$ , there is  $n_{\varepsilon,k} \in \mathbb{N}$  such that

$$\mu\left(\bigcup_{m=n_{\varepsilon,k}}^{\infty} \left\{|f_m - f| > \frac{1}{k}\right\}\right) < \frac{\varepsilon}{2^k}, \quad \forall m \geq n_{\varepsilon,k}.$$

Denote (the set of all possible divergence points! measurable!)

$$A_{\varepsilon} = \bigcup_{k=1}^{\infty} \bigcup_{m=n_{\varepsilon,k}}^{\infty} \left\{|f_m - f| > \frac{1}{k}\right\}.$$

Then  $\mu(A_\varepsilon) < \varepsilon$  and for any  $x \notin A_\varepsilon$ , we have  $\forall k$ ,

$$|f_m - f| \leq \frac{1}{k}, \quad \forall m > n_{\varepsilon, k}.$$

□

We have proved:

**Thm 19.** (1)

$$f_n \xrightarrow{a.u.} f \implies f_n \xrightarrow{a.e.} f \text{ and } f_n \xrightarrow{\mu} f$$

(2) *If  $\mu$  is finite, then*

$$f_n \xrightarrow{a.u.} f \iff f_n \xrightarrow{a.e.} f \implies f_n \xrightarrow{\mu} f$$

**Example 3.**

$$f_n(x) = \begin{cases} 1, & x \in (0, 1/n), \\ 0, & x \in [1/n, 1]. \end{cases}$$

**Example 4.**

$$f_n(x) = x^n, x \in [0, 1]$$

▷ 3. Let  $f = 0$  and  $f_n = 1_{A_n}$ . Then  $f_n \xrightarrow{\mu} f$  is equivalent to  $\mu(A_n) \rightarrow 0$  and  $\left(\lim_n f_n \neq f\right) = (A_n \text{ i.o.})$ .

Any sequence  $\{A_n\}$  so that  $\mu(A_n) \rightarrow 0$  but  $\mu(A_n \text{ i.o.}) > 0$  gives an example that  $f_n \xrightarrow{\mu} f \not\Rightarrow f_n \xrightarrow{a.e.} f$ . It is enough to have  $\mu(A_n) \rightarrow 0$  and

$$\sum_{i=1}^{\infty} 1_{A_n}(x) = \infty, \quad \sum_{i=1}^{\infty} 1_{A_n^c}(x) = \infty.$$

**Example 5.** For each  $n = 1, 2, \dots$  there is a unique decomposition  $n = k(k-1)/2 + i$  with  $k = 1, 2, \dots$ ,  $i = 1, 2, \dots, k$ .

$$f_n(x) = \begin{cases} 1, & x \in (((i-1)/k, i/k]), \\ 0, & \text{otherwise.} \end{cases}$$

**Example 6.** Consider

$$A_k^i = \left[ \frac{i-1}{k}, \frac{i}{k} \right], \quad h_k^i(x) = 1_{A_k^i}(x), \quad i = 1, \dots, k.$$

Let  $f_n$  be the sequence

$$\{h_1^1; h_2^1, h_2^2; h_3^1, h_3^2; h_3^3; \dots\}$$

**Thm 20.**  $f_n \xrightarrow{\mu} f \iff$  for any subsequence there is a further subsequence  $f_{n_k} \xrightarrow{a.u.} f$ .

PROOF. " $\implies$ ". Since any subsequence of  $f_n$  converges in measure to  $f$ , it is enough to show there is a subsequence  $f_{n_k} \xrightarrow{a.u.} f$ . To see this, for any  $k > 0$ , by definition of convergence in measure, we can choose  $n_k > n_{k-1}$  so that

$$\mu\left(|f_{n_k} - f| > \frac{1}{k}\right) \leq \frac{1}{2^k}.$$

Then

$$\mu\left(\bigcup_{k=m}^{\infty} |f_{n_k} - f| > \frac{1}{k}\right) \leq \sum_{k=m}^{\infty} \frac{1}{2^k} = \frac{1}{2^{m-1}}.$$

$\forall \varepsilon > 0$ , for large  $m$ ,

$$\bigcup_{k=m}^{\infty} \{|f_{n_k} - f| > \varepsilon\} \subset \bigcup_{k=m}^{\infty} \left\{ |f_{n_k} - f| > \frac{1}{k} \right\}.$$

So

$$\lim_m \mu \left( \bigcup_{k=m}^{\infty} |f_{n_k} - f| > \varepsilon \right) \leq \lim_m \mu \left( \bigcup_{k=m}^{\infty} |f_{n_k} - f| > \frac{1}{k} \right) = 0.$$

”  $\Leftarrow$  ” Suppose  $f_n \xrightarrow{\mu} f$  does not hold, i.e. there are  $n_k \rightarrow \infty$ ,  $\varepsilon_0 > 0$ ,  $\delta_0 > 0$  so that

$$\mu(|f_{n_k} - f| > \varepsilon_0) > \delta_0.$$

Then

$$\liminf_m \mu \left( \bigcup_{k=m}^{\infty} |f_{n_k} - f| > \varepsilon_0 \right) \geq \delta_0,$$

Contradicting Theorem [18](#).

□



Theorem 19 and Theorem 20 indicate that if  $f_n \xrightarrow{\mu} f$ , then there is a subsequence  $f_{n_k} \xrightarrow{a.e.} f$ .

## 4.2. 随机变量的分布函数.

**Def 19.**  $(\Omega, \mathcal{F}, P)$  is a probability space if  $P$  is a nonnegative measure on the  $\sigma$ -field  $\mathcal{F}$  with  $P(\Omega) = 1$ .

**Def 20.** A random variable (r.v.)  $X$  on  $(\Omega, \mathcal{F}, P)$  is a real-valued mapping,  $X : \omega \in \Omega \mapsto X(\omega) \in \mathbb{R}$ .

**Def 21.** The distribution function of a r.v.  $X$  is

$$F(x) = P(X \leq x).$$

Denoted by  $X \sim F$ .

**Thm 21.** Any distribution function  $F$  has the following properties.

(1) non-decreasing,  $F(-\infty) = 0$  and  $F(\infty) = 1$

(2) right continuity:  $\lim_{y \downarrow x} F(y) = F(x)$ .

(3) left limit exists:  $F(x-) = \lim_{y \uparrow x} F(y) = P(X < x)$ .

$$(4) \ P(X = x) = F(x) - F(x-).$$

The **inverse of the distribution function**  $F$  is defined as below.  
 $\forall z \in (0, 1)$ ,

$$(4.1) \qquad F^{-1}(z) = \inf\{x \in \mathbb{R} : F(x) \geq z\}.$$

▷ 4. *Also equivalently defined as,*

$$(4.2) \qquad F^{-1}(z) = \sup\{x \in \mathbb{R} : F(x) < z\}.$$

LEMMA 22.  $F^{-1}$  has the properties,

- (1)  $F^{-1}$  is real-valued non-decreasing.
- (2)  $F^{-1}$  is left-continuous and has right limit.
- (3)  $F^{-1}(F(x)) \leq x$ ,  $F(F^{-1}(z)) \geq z$ .
- (4)  $F^{-1}(z) \leq x$  iff  $F(x) \geq z$ .

PROOF. Exercise. □

**Thm 23.** *If  $F$  satisfies (1)(2)(3) of Theorem 21, there is a r.v.  $X$  with distribution  $F$ .*

PROOF. Let  $\Omega = (0, 1)$ ,  $\mathcal{F} = \mathcal{B}_{(0,1)}$  (i.e.  $(0, 1) \cap \mathcal{B}_{\mathbb{R}}$ ),  $P =$  Lebesgue measure. Define

$$X(\omega) = F^{-1}(\omega).$$

Then  $X$  is  $\mathcal{F}$ -measurable (check this!) and

$$\begin{aligned} P(\omega : X(\omega) \leq x) &= P(\omega : F(x) \geq \omega) \\ &= \text{Lebesgue measure of } (0, F(x)) = F(x). \end{aligned}$$

So  $X$  is a r.v. with distribution function  $F$ . □

▷ 5. *Another construction of a r.v.  $X$  with distribution  $F$  is to take  $(\Omega, \mathcal{F}) = (\mathbb{R}, \mathcal{B})$ ,  $P =$  the Lebesgue measure induced by  $F$  and consider the coordinate map  $X(\omega) = \omega$ .*

**4.3. 随机变量的收敛.** Probability space  $(\Omega, \mathcal{F}, P)$ , r.v.  $X_n, X$ ,

$$X_n \xrightarrow{a.s.} X \iff P(X_n = X) = 1.$$

$$X_n \xrightarrow{P} X \iff \forall \varepsilon > 0, \lim_n P(|X_n - X| > \varepsilon) = 0.$$

**Def 22.**  $X_n \sim F_n, X \sim F$ . Convergence in distribution (weak convergence):  $F_n(x) \rightarrow F(x)$  for all  $x$  where  $F$  is continuous, written  $X_n \xrightarrow{d} X$ .

**Thm 24.**  $X_n \sim F_n, X \sim F$ .

$$X_n \xrightarrow{a.s.} X \implies X_n \xrightarrow{P} X \implies X_n \xrightarrow{d} X.$$

PROOF. 1. The first implication is a special case of Theorem 17.

2. Check the second implication.  $\forall \varepsilon, x \in \mathbb{R}, n \in \mathbb{N}$ ,

$$\begin{aligned} & P(X \leq x - \varepsilon) - P(|X_n - X| > \varepsilon) \\ & \leq P(X_n \leq x) \\ & \leq P(X_n \leq x, |X_n - X| \leq \varepsilon) + P(X_n \leq x, |X_n - X| > \varepsilon) \\ & \leq P(X \leq x + \varepsilon) + P(|X_n - X| > \varepsilon). \end{aligned}$$

So  $n \rightarrow \infty, \varepsilon \rightarrow 0$  yield

$$F(x-) \leq \liminf_n P(X_n \leq x) \leq \limsup_n P(X_n \leq x) \leq F(x).$$

□

LEMMA 25.  $F_n \xrightarrow{w} F \iff F_n^{-1} \xrightarrow{w} F^{-1}$ .

PROOF OF "  $\implies$  ". Construct r.v.s'  $X_n \sim F_n$ ,  $X \sim F$  as Theorem 23. Fix any  $\omega$ .

1. Choose any  $\varepsilon > 0$  so that  $F$  is continuous at  $X(\omega) - \varepsilon$  (the discontinuities of  $F$  are at most countable,  $\varepsilon$  can be arbitrarily small). By the definition (the infimum!) of  $X(\omega)$ ,

$$F(X(\omega) - \varepsilon) < \omega.$$

Then, for large  $n$ ,

$$F_n(X(\omega) - \varepsilon) < \omega.$$

so (note the above inequality is strict)

$$X(\omega) - \varepsilon \leq X_n(\omega).$$

Hence

$$X(\omega) \leq \liminf_n X_n(\omega).$$

2. To see the opposite. Choose any  $\varepsilon, \delta > 0$  so that  $X$  is continuous at  $\omega$  and  $F$  is continuous at  $X(\omega) + \varepsilon$ , then by Lemma 22

$$F(X(\omega + \delta) + \varepsilon) \geq F(X(\omega + \delta)) \geq \omega + \delta > \omega.$$

For large  $n$  ( $\delta > 0$ ),

$$F_n(X(\omega + \delta) + \varepsilon) \geq \omega.$$

By Lemma 22 again,

$$X(\omega + \delta) + \varepsilon \geq X_n(F_n(X(\omega + \delta) + \varepsilon)) \geq X_n(\omega).$$

Let  $n \rightarrow \infty$ ,  $\varepsilon \rightarrow 0$ ,  $\delta \rightarrow 0$  (continuity at  $\omega$ ),

$$X(\omega) \geq \limsup_n X_n(\omega).$$

□

**Thm 26** (Skorohod).  $X_n \sim F_n$ ,  $X \sim F$ . Suppose  $X_n \xrightarrow{d} X$ . There exist r.v.  $\bar{X}_n, \bar{X}$  on a common probability space so that  $\bar{X}_n \stackrel{d}{=} X_n$ ,  $\bar{X} \stackrel{d}{=} X$ ,  $\bar{X}_n \xrightarrow{a.s.} \bar{X}$ .

PROOF. Let  $\Omega = (0, 1)$ ,  $\mathcal{F} = \mathcal{B}_{(0,1)}$ ,  $P =$  Lebesgue measure. By Theorem 23 there exist r.v. on  $(\Omega, \mathcal{F}, P)$  so that  $\bar{X}_n \sim F_n$ ,  $\bar{X} \sim F$ . Lemma 25 then says  $F_n^{-1} \xrightarrow{w} F^{-1}$ . Since the discontinuity set of  $F^{-1}$  is countable,  $F_n^{-1}(\omega) \rightarrow F^{-1}(\omega)$  for almost all  $\omega \in \Omega$ , i.e.  $\bar{X}_n(\omega) \xrightarrow{a.s.} \bar{X}(\omega)$ .  $\square$

## 5. 积分

**5.1. 非负可测函数积分.**  $(E, \mathcal{F}, \mu)$  a measure space,  $f \in \mathcal{F}$  with values in  $[0, \infty]$ ,. A *finite (measurable) partition* of  $E$  is a finite collection of  $\mathcal{F}$ -measurable sets  $\{A_i : i = 1, \dots, m\}$  with  $\bigcup_{i=1}^m A_i = E$ .

$$(5.1) \quad \int f d\mu \triangleq \sup_{\text{finite partitions}} \sum_i \left[ \inf_{x \in A_i} f(x) \right] \mu(A_i).$$

Convention:  $0 \cdot \infty = 0$ .

▷ 6. Consider

$$(5.2) \quad \int f d\mu \triangleq \inf_{\text{finite partitions}} \sum_i \left[ \sup_{x \in A_i} f(x) \right] \mu(A_i).$$

Is (5.2) a good definition of integration?

**Properties:**  $f, g \in \mathcal{F}$  nonnegative.

(1) If  $f = 0$ ,  $\mu$ -a.e., then  $\int f d\mu = 0$ .



(2) If  $\mu(f > 0) > 0$ , then  $\int f d\mu > 0$ .

(3) If  $\int f d\mu < \infty$ , then  $f < \infty, \mu$ -a.e.

(4) If  $f \leq g, \mu$ -a.e., then  $\int f d\mu \leq \int g d\mu$ .

(5) If  $f = g, \mu$ -a.e., then  $\int f d\mu = \int g d\mu$ .

**Thm 27** (Monotone convergence Theorem). *If  $0 \leq f_n \uparrow f, \mu$ -a.e., then  $0 \leq \int f_n d\mu \uparrow \int f d\mu$ .*

PROOF. 1. First prove it under the assumption that

$$0 \leq f_n(x) \uparrow f(x), \forall x.$$

Integration is monotonic, so  $\int f_n d\mu \leq \int f d\mu$ . It remains to show

$$(5.3) \quad \lim_n \int f_n d\mu \geq \int f d\mu$$

or

$$\lim_n \int f_n d\mu \geq S = \sum_{i=1}^m c_i \mu(A_i)$$

for any finite measurable partition  $\{A_i : i = 1, \dots, m\}$  and  $c_i = \inf_{A_i} f$ .

For such a partition, assume that the sum  $S$ ,  $c_i$  and  $\mu(A_i)$  are all finite. Fix  $\alpha < 1$ , define

$$A_{i,n} = \{x \in A_i : f_n(x) > \alpha c_i\}.$$

Since  $f_n \uparrow f$ ,  $A_{i,n} \uparrow A_i$ . Consider the *measurable* partition

$$\{A_{i,n} : i = 1, \dots, m\} \cup \left\{ \left( \bigcup_{i=1}^m A_{i,n} \right)^c \right\}.$$

Then

$$\int f_n d\mu \geq \sum_{i=1}^m \alpha c_i \mu(A_{i,n}).$$

Let  $n \rightarrow \infty$  and use continuity from below,

$$\lim_n \int f_n d\mu \geq \sum_{i=1}^m \alpha c_i \mu(A_i).$$

Finally let  $\alpha \rightarrow 1$ , (5.3) is proved.

Now suppose  $S$  is finite but not all of  $c_i, \mu(A_i)$ . Then  $c_i \mu(A_i)$ ,  $i = 1, \dots, m$  are finite.  $c_i$  or  $\mu(A_i)$  may be infinity, but then  $c_i \mu(A_i)$  must be zero. Use the adjusted partition  $\{A_i : c_i \mu(A_i) > 0\} \cup \{\text{complement}\}$ .

Lastly suppose  $S$  is infinite. Then there is some  $i_0$ ,  $c_{i_0} \mu(A_{i_0}) = \infty$ , i.e.,  $c_{i_0} > 0$ ,  $\mu(A_{i_0}) > 0$  and at least one of them is  $\infty$ . In this case

$$\int f d\mu = \infty.$$

To prove (5.3), let  $a, b$  satisfy

$$0 < a < c_{i_0} \leq \infty, \quad 0 < b < \mu(A_{i_0}) \leq \infty.$$

Define

$$A_{i_0, n} = \{x \in A_{i_0} : f_n(x) > a\}.$$

Since  $f_n \uparrow f$ ,  $A_{i_0,n} \uparrow A_{i_0}$  and  $\mu(A_{i_0,n}) > b$  for  $n$  larger than some  $n_{a,b}$ . For the partition  $\{A_{i_0,n}, A_{i_0,n}^c\}$ , we have

$$\int f_n d\mu \geq a\mu(A_{i_0,n}) > ab, \forall n > n_{a,b}.$$

Let  $a \rightarrow \infty$  if  $c_{i_0} = \infty$ ,  $b \rightarrow \infty$  if  $\mu(A_{i_0,n}) = \infty$ , we get

$$\lim_n \int f_n d\mu = \infty.$$

**2.** If  $0 \leq f_n \uparrow f$  on  $A$  with  $\mu(A^c) = 0$ , then  $0 \leq f_n 1_A \uparrow f 1_A$  holds everywhere. Then apply step **1**.  $\square$

**5.2. 可测函数积分.**  $f \in \mathcal{F}$  with values in  $[-\infty, \infty]$ ,

$$\int f d\mu \triangleq \int f^+ d\mu - \int f^- d\mu.$$

$f$  is said to be integrable if  $\int f^+ d\mu, \int f^- d\mu$  are finite. So  $f$  integrable iff  $|f|$  integrable.

**Properties:**  $f, g \in \mathcal{F}$  integrable.

(1) If  $f \leq g$ ,  $\mu$ -a.e., then  $\int f d\mu \leq \int g d\mu$ .

(2) If  $\alpha, \beta \in \mathbb{R}$ , then  $\alpha f + \beta g$  is integrable,

$$\int (\alpha f + \beta g) d\mu = \alpha \int f d\mu + \beta \int g d\mu.$$

**Example 7.** Let  $E = \{1, 2, 3, \dots\}$ ,  $\mathcal{F} = \{\text{all subsets of } E\}$ ,  $\mu =$  counting measure. A function on  $E$  is a sequence  $x_1, x_2, \dots$ . Any function is  $\mathcal{F}$ -measurable.  $\{x_k : k = 1, 2, \dots\}$  is  $\mu$ -integrable if and only if  $\sum_{k=1}^{\infty} |x_k|$  converges. When  $\mu$ -integrable,

$$\sum_{k=1}^{\infty} |x_k| = \sum_{k=1}^{\infty} x_k^+ - \sum_{k=1}^{\infty} x_k^-.$$

The function  $x_k = (-1)^{k+1}/k$ ,  $k = 1, 2, \dots$  is not  $\mu$ -integrable, although

$$\lim_m \sum_{k=1}^m (-1)^{k+1} \frac{1}{k} = \ln 2.$$

**Thm 28** (Fatou's lemma). *Given  $f_n$  measurable.*

(1) *If  $g$  integrable,  $f_n \geq g$ ,  $\mu$ -a.e, then  $\liminf_n f_n$  is integrable and*

$$\int \liminf_n f_n d\mu \leq \liminf_n \int f_n d\mu.$$

(1) *If  $g$  integrable,  $f_n \leq g$ ,  $\mu$ -a.e, then  $\limsup_n f_n$  is integrable and*

$$\limsup_n \int f_n d\mu \leq \int \limsup_n f_n d\mu.$$

**Thm 29** (Lebesgue's dominated convergence theorem). *Given  $g$  nonnegative integrable,  $|f_n| \leq g$ ,  $\mu$ -a.e.. If  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ , then*

$$\int f_n d\mu \longrightarrow \int f d\mu.$$

The following is a generalized dominated convergence theorem.

**Thm 30.** Given  $g_n$  nonnegative integrable,  $|f_n| \leq g_n$ ,  $\mu$ -a.e. with  $g_n \xrightarrow{a.e.} g$  and  $\int g_n d\mu \longrightarrow \int g d\mu$ . If  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ , then

$$\int f_n d\mu \longrightarrow \int f d\mu.$$

**Example 8** (Weierstrass M-test). If  $|x_{n,m}| \leq M_m$ ,  $\sum_{m=1}^{\infty} M_m < \infty$ ,

$\lim_n x_{n,m} = x_m$  for each  $m$ . Then

$$\lim_n \sum_{m=1}^{\infty} x_{n,m} = \sum_{m=1}^{\infty} x_m.$$

**Example 9** (Bounded convergence theorem). Suppose  $\mu$  is finite,  $M > 0$ .  $|f_n| \leq M$ ,  $\mu$ -a.e.. If  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ , then

$$\int f_n d\mu \longrightarrow \int f d\mu.$$

**Example 10.** If  $f_n \geq 0$  or  $\sum_{n=1}^{\infty} \int |f_n| d\mu < \infty$ , then

$$\int \sum_{n=1}^{\infty} f_n d\mu = \sum_{n=1}^{\infty} \int f_n d\mu.$$

From this we get

**Example 11.** If  $x_{n,m} \geq 0$  or  $\sum_{m=1}^{\infty} \sum_{n=1}^{\infty} |x_{n,m}| < \infty$ , then

$$\sum_{m=1}^{\infty} \sum_{n=1}^{\infty} x_{n,m} = \sum_{n=1}^{\infty} \sum_{m=1}^{\infty} x_{n,m}.$$

**Example 12** (Abel's theorem). Suppose that the series  $\sum_{k=1}^{\infty} |c_k|$  is convergent. Then

$$\lim_{x \rightarrow 1^-} \sum_{k=1}^{\infty} c_k x^k = \sum_{k=1}^{\infty} c_k.$$



**5.3. Change of variables.**  $(E_1, \mathcal{F}_1)$ ,  $(E_2, \mathcal{F}_2)$  are measurable spaces,  $\mu$  is a measure on  $\mathcal{F}_1$ .  $T$  is measurable mapping from  $(E_1, \mathcal{F}_1)$  to  $(E_2, \mathcal{F}_2)$ . Define

$$(5.4) \quad \nu(B) = \mu(T^{-1}(B)), \quad \forall B \in \mathcal{F}_2.$$

Then  $\nu(B)$  is a measure on  $\mathcal{F}_2$  and for any  $f \in \mathcal{F}_2$ ,

$$(5.5) \quad \int_{E_2} f d\nu = \int_{E_1} f \circ T d\mu.$$

Note if  $f = 1_B$ , then  $f \circ T(x) = 1_B(T(x)) = 1_{T^{-1}(B)}(x)$ , since  $T(x) \in B$  iff  $x \in T^{-1}(B)$ . So in this case (5.5) reduces to (5.4).

## 6. $L_p$ 空间

### 6.1. Inequalities.

LEMMA 31 (Jensen's inequality). *Let  $(\Omega, \mathcal{F}, \mu)$  be a measure space with  $\mu(\Omega) = 1$ ,  $X$  a  $\mu$ -integrable function on  $\Omega$ ,  $\varphi$  convex on  $\mathbb{R}$ . Then*

$$(6.1) \quad \varphi\left(\int_{\Omega} X d\mu\right) \leq \int_{\Omega} \varphi(X) d\mu.$$

Equality holds iff  $\varphi$  is linear on some convex set  $A \subset \mathbb{R}$  with  $\mu(X^{-1}A) = 1$ .

PROOF. Denote by  $\mu_X$  the induced measure of  $X$  on  $\mathbb{R}$  (ref section 5.3), then (6.1) is equivalent to

$$(6.2) \quad \varphi\left(\int_{\mathbb{R}} x d\mu_X\right) \leq \int_{\mathbb{R}} \varphi(x) d\mu_X$$

(Apply (5.5) with  $f(x) = x$ ,  $T = X$ ). It is enough to prove (6.2).

1. Denote  $\bar{x} = \int_{\mathbb{R}} x d\mu_X$ . Since  $\varphi$  is convex, there is a supporting line  $L(x) = ax + b$  through  $\bar{x}$ , i.e.  $L(\bar{x}) = \varphi(\bar{x})$  and

$$L(x) \leq \varphi(x), \quad \forall x.$$

Then

$$(6.3) \quad \int_{\mathbb{R}} L(x) d\mu_X \leq \int_{\mathbb{R}} \varphi(x) d\mu_X.$$

The LHS equals  $\varphi\left(\int_{\mathbb{R}} x d\mu_X\right)$ , hence (6.2) follows.

2. Suppose the equality in (6.2) holds, then by the above computation

$$\int_{\mathbb{R}} [\varphi(x) - L(x)] d\mu_X = 0.$$

The integrand is nonnegative, so the measurable set

$$A = \{x \in \mathbb{R} : \varphi(x) - L(x) = 0\}$$

has full measure, i.e.  $\mu_X(A) = 1$ . Moreover the set  $A$  is convex (verify directly!). On the other hand, if  $\varphi$  is linear on some convex  $A \subset \mathbb{R}$  with  $\mu(X^{-1}A) = 1$ , then  $\mu_X(A) = 1$ ,

$$\int_{\mathbb{R}} L(x) d\mu_X = \int_A L(x) d\mu_X, \quad \int_{\mathbb{R}} \varphi(x) d\mu_X = \int_A \varphi(x) d\mu_X.$$

Hence by (6.3),

$$\int_A [\varphi(X) - L(X)] d\mu \geq 0.$$

But the integrand  $\varphi - L$  is nonnegative and linear on  $A$ . Since  $A \subset \mathbb{R}$  is convex, it must be an interval. So the above integral is zero, hence the equality of (6.2) holds.  $\square$

Notice that Lemma 31 does not require  $\varphi(X)$  to be  $\mu$ -integrable. From (6.3) it is clear that either  $\int_{\Omega} \varphi(X) d\mu$  exists or equals infinity, in the latter case (6.1) trivially holds.

LEMMA 32.  $a, b \in \mathbb{R}$ ,  $1 \leq p < \infty$ ,

$$|a + b|^p \leq 2^{p-1}(|a|^p + |b|^p).$$

PROOF. Apply Jensen's inequality with  $\varphi(x) = |x|^p$ ,

$$\left| \frac{a + b}{2} \right|^p \leq \frac{|a|^p + |b|^p}{2}.$$

□

LEMMA 33 (Young's inequality).  $a, b \geq 0$ ,  $1 < p, q < \infty$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ ,

$$a^{1/p} b^{1/q} \leq \frac{a}{p} + \frac{b}{q}.$$

*Equal iff  $a = b$ .*

PROOF. The inequality holds if  $ab = 0$ . In this case equality holds iff  $a = b = 0$ . Now suppose  $ab > 0$ . Apply Jensen's inequality with  $\varphi(x) = -\ln x$ ,

$$-\ln\left(\frac{a}{p} + \frac{b}{q}\right) \leq -\frac{1}{p}\ln a - \frac{1}{q}\ln b.$$

Since  $\varphi$  is strictly convex (can touch a linear function at exactly one point), equality holds iff  $a = b$ .  $\square$

$(E, \mathcal{F}, \mu)$  is a measure space in the following definitions.

**Def 23.**  $p = 1$ , let

$$L_1 \triangleq \{f \in \mathcal{F} : |f| \text{ is } \mu\text{-integrable}\}$$

and

$$\|f\|_1 = \|f\|_{L_1} = \int |f| d\mu.$$

**Def 24.**  $1 < p < \infty$ , let

$$L_p \triangleq \{f \in \mathcal{F} : |f|^p \in L_1\}$$

and

$$\|f\|_p = \|f\|_{L_p} = \left( \int |f|^p d\mu \right)^{1/p}.$$

**Def 25.**  $p = \infty$ , let

$$L_\infty \triangleq \{f \in \mathcal{F} : \text{there is } C > 0 \text{ such that } |f| \leq C, \text{ a.e.}\}$$

and

$$\|f\|_\infty = \|f\|_{L_\infty} = \inf\{C : |f| \leq C, \text{ a.e.}\}.$$

We could have written  $L_p(\mu)$  to emphasize the dependence of the spaces  $L_p$  on the measure  $\mu$ . But, when no ambiguity arises from the contexts, we will simply drop  $\mu$  from the notation.

**Thm 34** (Hölder inequality).  $1 \leq p, q \leq \infty$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ ,  $f \in L_p$ ,  $g \in L_q$ , then  $fg \in L_1$  and

$$(6.4) \quad \|fg\|_1 \leq \|f\|_p \|g\|_q.$$

If  $p = 1$ , equality iff  $|g| = \|g\|_\infty$ , a.e. on the set where  $f \neq 0$ .

If  $p = \infty$ , equality iff  $|f| = \|f\|_\infty$ , a.e. on the set where  $g \neq 0$ .

If  $1 < p < \infty$ , equality iff there are nonnegative constants  $\alpha, \beta$  such that  $(\alpha, \beta) \neq (0, 0)$ ,  $\alpha|f|^p = \beta|g|^q$ , a.e.

PROOF. **1.** The inequality easily follows if  $p = 1$  or  $p = \infty$ . To see the equality, suppose  $p = 1$ , then  $q = \infty$ . (6.4) is equivalent to

$$\int |f|(\|g\|_\infty - |g|) \geq 0.$$

It is equality iff  $|g| = \|g\|_\infty$ , a.e. on the set where  $f \neq 0$ .

**2.** Suppose  $1 < p, q < \infty$ . The conclusion is obvious if  $\|f\|_p = 0$  or  $\|g\|_q = 0$ . Hence we assume that  $0 < \|f\|_p, \|g\|_q < \infty$ . Using Young's inequality with

$$a = \left( \frac{|f|}{\|f\|_p} \right)^p, \quad b = \left( \frac{|g|}{\|g\|_q} \right)^q,$$

we have

$$\frac{|fg|}{\|f\|_p\|g\|_q} \leq \frac{1}{p} \left( \frac{|f|}{\|f\|_p} \right)^p + \frac{1}{q} \left( \frac{|g|}{\|g\|_q} \right)^q, \text{ a.e.}$$

Integrating on both sides gives

$$\int \frac{|fg|}{\|f\|_p\|g\|_q} d\mu \leq \frac{1}{p} + \frac{1}{q} = 1,$$

which is the desired inequality. The equality holds iff  $a = b$ , *a.e.* i.e.,

$$\|g\|_q^q |f|^p = \|f\|_p^p |g|^q, \text{ a.e.}$$

□

A familiar case of Hölder inequality is the following.

**Thm 35** (Cauchy–Schwarz inequality).  *$f, g \in L_2$ , then  $fg \in L_1$  and*

$$\|fg\|_1 \leq \|f\|_2 \|g\|_2.$$



**Thm 36** (Minkowski inequality).  $1 \leq p \leq \infty$ ,  $f, g \in L_p$ , then  $f + g \in L_p$  and

$$(6.5) \quad \|f + g\|_p \leq \|f\|_p + \|g\|_p.$$

*If  $p = 1$  or  $p = \infty$ , equality iff  $fg \geq 0$ , a.e..*

*If  $1 < p < \infty$ , equality iff there are nonnegative constants  $\alpha, \beta$  such that  $(\alpha, \beta) \neq (0, 0)$ ,  $\alpha f = \beta g$ , a.e.*

**PROOF. 1.** The case  $p = 1$  or  $p = \infty$  is immediate.

**2.** Suppose  $1 < p < \infty$ . Let  $q > 1$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ . By Hölder inequality

$$\begin{aligned} \|f + g\|_p^p &= \int |f + g| |f + g|^{p-1} \leq_{(e1)} \int |f| |f + g|^{p-1} + \int |g| |f + g|^{p-1} \\ &\leq_{(e2)} \|f\|_p \| |f + g|^{p-1} \|_q + \|g\|_p \| |f + g|^{p-1} \|_q \\ &= \|f\|_p \|f + g\|_p^{p-1} + \|g\|_p \|f + g\|_p^{p-1} \end{aligned}$$

Here

$$\begin{aligned}\| |f + g|^{p-1} \|_q &= \left( \int (|f + g|^{p-1})^q \right)^{1/q} = \left( \int |f + g|^p \right)^{1/q} \\ &= \|f + g\|_p^{p/q} = \|f + g\|_p^{p-1}.\end{aligned}$$

(e1) is equality iff  $fg \geq 0$ , *a.e.*, (e2) is equality iff there are nonnegative constants  $a, b, c, d$  such that  $(a, b) \neq (0, 0)$ ,  $(c, d) \neq (0, 0)$ ,

$$a|f|^p = b(|f + g|^{p-1})^q, \quad c|g|^p = d(|f + g|^{p-1})^q, \quad \text{a.e.}$$

Hence

$$a|f| = b|f + g|, \quad c|g| = d|f + g|, \quad \text{a.e.}$$

The conclusion follows by combining the equality conditions of (e1)(e2).  $\square$

**Def 26.**  $0 < p < 1$ , let

$$L_p \triangleq \left\{ f \in \mathcal{F} : \int |f|^p d\mu < \infty \right\}$$

and

$$\|f\|_p = \int |f|^p d\mu.$$

LEMMA 37. Let  $a, b \in \mathbb{R}$ ,  $0 < p < 1$ .  $|a + b|^p \leq |a|^p + |b|^p$ .

PROOF. Since  $||a| + |b||^p \leq |a|^p + |b|^p$  implies the desired inequality, we assume w.l.g. that  $a, b$  are of the same sign. Suppose  $a \neq 0$ , otherwise there is nothing to prove. Finally it suffices to show that

$$(1 + s)^p \leq 1 + s^p, \quad s \geq 0,$$

which is verified by elementary calculus. □

Lemma 32 and Lemma 37 can be merged into the compact form,

$$(6.6) \quad |a + b|^p \leq C_p(|a|^p + |b|^p), \quad 0 < p < \infty,$$

where  $C_p = 2^{p-1} \vee 1$ .

**Thm 38.**  $0 < p < 1$ ,  $\|f + g\|_p \leq \|f\|_p + \|g\|_p$ .

## 6.2. Completeness.

**Thm 39.** *Let  $0 < p \leq \infty$ ,  $L_p$  is complete.*

PROOF FOR  $p = \infty$ . Let  $f_n \in L_\infty$ . Suppose that  $f_n$  is Cauchy. Given  $k \geq 1$ , there is  $n_k$  such that

$$\|f_m - f_n\|_\infty \leq \frac{1}{k}, \quad \forall m, n > n_k.$$

Hence there is a null set<sup>1</sup>  $A_k$  such that

$$|f_m - f_n| \leq \frac{1}{k}, \quad \forall x \in A_k^c, \quad m, n > n_k.$$

Then  $A = \bigcup_{k=1}^{\infty} A_k$  is a null set and  $f_n(x)$  is Cauchy for each  $x \in A^c$ .

Hence there exist  $f$ ,  $f_n \rightarrow f$  for  $x \in A^c$ . Let  $m \rightarrow \infty$  in the above inequality we get

$$|f_n - f| \leq \frac{1}{k}, \quad \forall x \in A^c, \quad n > n_k.$$

---

<sup>1</sup>A null set is a measurable set with measure zero.

So  $f \in L_\infty$  and

$$\|f_n - f\|_\infty \leq \frac{1}{k}, \quad \forall n > n_k.$$

Therefore  $f_n$  converges to  $f$  in  $L_\infty$ . □

PROOF FOR  $0 < p < \infty$ . Let  $f_n \in L_p$ . Suppose that  $f_n$  is Cauchy in  $L_p$ ,

$$(6.7) \quad \lim_{m, n \rightarrow \infty} \|f_m - f_n\|_p = 0.$$

We intend to show that  $\lim_{n \rightarrow \infty} \|f_n - f\|_p = 0$  for some  $f \in L_p$ . Owing to (6.7), we have a subsequence  $n_k \rightarrow \infty$  so that

$$(6.8) \quad \|f_{n_{k+1}} - f_{n_k}\|_p < \frac{1}{2^k}.$$

We claim that

- (a) there is  $h \in L_p$  such that  $|f_{n_k}| \leq h$ , *a.e.*
- (b)  $\lim_k f_{n_k} \rightarrow f$ , *a.e.* for some  $f \in L_p$ .
- (c)  $\lim_k \|f_{n_k} - f\|_p = 0$ .

The conclusion of the Theorem clearly follows once (c) is proved, since a Cauchy sequence converges iff it has a convergent subsequence. Let

$$g_k = \sum_{i=1}^k |f_{n_{i+1}} - f_{n_i}|, \quad g = \sum_{i=1}^{\infty} |f_{n_{i+1}} - f_{n_i}|.$$

Then  $0 \leq g_k \uparrow g$  and  $\|g_k\|_p \leq 1$  by (6.8) (Theorem 36 or Theorem 38). Using monotone convergence theorem,

$$\int g^p d\mu = \lim_k \int (g_k)^p d\mu \leq 1.$$

This shows  $g \in L_p$  and that  $g < \infty$ , *a.e.* Therefore

$$f_{n_k} = f_{n_1} + \sum_{i=1}^k (f_{n_{i+1}} - f_{n_i})$$

converges almost everywhere to some measurable function  $f$  and

$$|f_{n_k}| \leq |f_{n_1}| + g.$$

Let  $k \rightarrow \infty$ , we have

$$|f| \leq |f_{n_1}| + g, \text{ a.e.}$$

hence  $f \in L_p$ . (a)(b) follows with  $h = |f_{n_1}| + g$ . By inequality (6.6),

$$\begin{aligned} |f_{n_k} - f|^p &\leq C_p(|f_{n_k}|^p + |f|^p) \leq C_p(|f_{n_1}| + g)^p + |f|^p \\ &\leq C_p(C_p(|f_{n_1}|^p + g^p) + |f|^p). \end{aligned}$$

Therefore (c) is a result of the dominated convergence theorem. □

**COROLLARY 1.** (1)  $0 < p < 1$ ,  $L_p$  is a complete metric space.  
 (2)  $1 \leq p \leq \infty$ ,  $L_p$  is a Banach space.

### 6.3. $L_p$ and weak convergence.

**Thm 40.** Let  $0 < p < \infty$ ,  $f_n \in L_p$ ,  $f \in L_p$ .

- (1)  $f_n \xrightarrow{L_p} f \implies f_n \xrightarrow{\mu} f$  and  $\|f_n\|_p \rightarrow \|f\|_p$ .  
 (2)  $f_n \xrightarrow{\text{a.e.}} f$  or  $f_n \xrightarrow{\mu} f$ , then

$$\|f_n\|_p \rightarrow \|f\|_p \iff f_n \xrightarrow{L_p} f.$$

PROOF. **1.** To prove (1), use Markov inequality

$$\mu(|f_n - f| > \varepsilon) \leq \frac{1}{\varepsilon^p} \|f_n - f\|_p^p.$$

and the triangle inequality

$$\left| \|f_n\|_p - \|f\|_p \right| \leq \|f_n - f\|_p.$$

**2.** "  $\Leftarrow$  " of (2) is included in step **1**.

**3.** "  $\Rightarrow$  " of (2). In view of Theorem [20](#), it is enough to prove the case where  $f_n \xrightarrow{a.e.} f$ . Define

$$g_n = C_p(|f_n|^p + |f|^p) - |f_n - f|^p,$$



where  $C_p = 2^{p-1} \vee 1$ . Then  $g_n \geq 0$  by inequality (6.6) and  $\lim_n g_n = 2C_p|f|^p$ , a.e. Using Fatou's lemma

$$\begin{aligned} \int 2C_p|f|^p d\mu &= \int \lim_n g_n d\mu \leq \liminf_n \int g_n d\mu \\ &= \int 2C_p|f|^p d\mu - \limsup_n \int |f_n - f|^p. \end{aligned}$$

Canceling  $\int 2C_p|f|^p d\mu$  from both side gives

$$\lim_n \int |f_n - f|^p = 0.$$

□

**Def 27.**  $(E, \mathcal{F}, \mu)$  is a measure space.  $1 \leq p < \infty$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ ,  $f_n$  converges weakly to  $f$  in  $L_p$ , denoted by  $f_n \xrightarrow{w-L_p} f$ , if

$$\lim_n \int f_n g d\mu = \int f g d\mu, \quad \forall g \in L_q.$$

$\mu$  is additionally assumed to be  $\sigma$ -finite if  $p = 1$ .

**Thm 41.**  $1 \leq p < \infty$ .  $f_n \xrightarrow{L_p} f$  implies  $f_n \xrightarrow{w-L_p} f$ .

PROOF. By Hölder inequality (Theorem 34),  $\forall g \in L_q$ ,  $q$  conjugate to  $p$ ,

$$\int |f_n - f| |g| d\mu \leq \|f_n - f\|_p \|g\|_q.$$

□

**Thm 42.**  $(E, \mathcal{F}, \mu)$  is a measure space. Let  $1 < p < \infty$ ,  $\{f_n\}$  bounded in  $L_p$ . If  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$  for some measurable  $f$ , then  $f \in L_p$  and  $f_n \xrightarrow{w-L_p} f$ .

PROOF. Let  $g \in L_q$ ,  $q$  conjugate to  $p$ . As before, it is enough to prove it for  $f_n \xrightarrow{a.e.} f$ .

1.  $f \in L_p$  is a consequence of Fatou's lemma,

$$\int |f|^p d\mu = \int \lim_n |f_n|^p d\mu \leq \liminf_n \int |f_n|^p d\mu \leq \sup_n \|f_n\|_{L_p}^p < \infty.$$

It follows that  $\{f_n - f\}$  is bounded in  $L_p$ .

**2.** Fix  $\varepsilon > 0$ , let  $\delta > 0$ , define  $A_\delta = \{x \in E : \delta \leq |g|^q \leq 1/\delta\}$  and write

$$\int |f_n - f| |g| d\mu = \int_{A_\delta \cap B} + \int_{A_\delta \cap B^c} + \int_{A_\delta^c}.$$

Choose  $\delta$  small so that

$$\int_{A_\delta^c} \leq \|f_n - f\|_p \|g 1_{A_\delta^c}\|_q < \frac{\varepsilon}{3}.$$

With  $\delta$  fixed, we have

$$\int_{A_\delta \cap B^c} \leq \|f_n - f\|_p \|g 1_{A_\delta \cap B^c}\|_q < \frac{\varepsilon}{3},$$

as soon as  $B \subset A_\delta$  is such that  $\mu(A_\delta \cap B^c)$  is smaller than some  $\varepsilon'$ .

Note  $|g| \leq 1/\delta^{1/q}$  on  $A_\delta$ . Since  $\mu(A_\delta)$  is finite by Markov inequality, so a subset  $B \subset A_\delta$  can be chosen so that  $\mu(A_\delta \cap B^c) < \varepsilon'$  and  $|f_n - f|$

converges uniformly to 0 on  $A_\delta \cap B$  (Theorem 19). Hence for large  $n$ ,

$$\int_{A_\delta \cap B} \leq \frac{1}{\delta^{1/q}} \int_{A_\delta \cap B} |f_n - f| d\mu < \frac{\varepsilon}{3}.$$

□

Note the above proof does not get through if  $p = 1$  (so that  $q = \infty$ ). The example below demonstrates, in general, Theorem 42 does not for  $p = 1$ .

**Example 13.**  $E = (0, 1)$  with the usual Lebesgue measure,  $f_n = n1_{(0, 1/n)}$ . Clearly  $\|f_n\|_1 = 1$ ,  $f_n \xrightarrow{\mu} f = 0$ . But with  $g = 1 \in L_\infty$ ,  $\lim_n \int f_n g d\mu = 1 \neq 0 = \int f g d\mu$ , hence  $f_n \xrightarrow{w-L_1} f$  does not hold.

However we have

**Thm 43.**  $(E, \mathcal{F}, \mu)$  is a measure space. Let  $\{f_n\} \in L_1$ . Suppose  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ . Then

$$f \in L_1, \quad \|f_n\|_1 \rightarrow \|f\|_1 \iff f_n \xrightarrow{L_1} f.$$

Either of them gives  $\int_A f_n d\mu \rightarrow \int_A f d\mu, \forall A \in \mathcal{F}$ .

PROOF. The first conclusion is contained in Theorem 40. So  $f_n \xrightarrow{w-L^2} f$  by Theorem 41. To complete the proof, take  $1_A \in L_\infty$  as test function.  $\square$

**6.4. Uniform integrability.** Let  $(E, \mathcal{F}, \mu)$  be a measure space.

**Def 28.**  $\mathcal{H} = \{f_t : t \in T\}$  is uniformly integrable if

$$(6.9) \quad \lim_{a \rightarrow \infty} \sup_{f \in \mathcal{H}} \int_{\{|f| \geq a\}} |f| d\mu = 0.$$

**Def 29.**  $\mathcal{H} = \{f_t : t \in T\}$  is absolutely continuous if,  $\forall \varepsilon > 0$ , there is  $\delta > 0$  so that

$$\sup_{f \in \mathcal{H}} \int 1_A |f| d\mu < \varepsilon \text{ for any } A \text{ with } \mu(A) < \delta.$$

**Thm 44.** Suppose  $(E, \mathcal{F}, \mu)$  is a measure space with  $\mu$  finite.  $\mathcal{H} = \{f_t : t \in T\}$  is uniformly integrable if and only if  $\mathcal{H}$  is absolutely continuous and bounded in  $L_1$ .

PROOF. 1. If  $\mathcal{H}$  is uniformly integrable,  $\forall \varepsilon > 0$ , there is  $a_0 > 0$  so that

$$\sup_{f \in \mathcal{H}} \int_{\{|f| \geq a\}} |f| d\mu \leq \frac{\varepsilon}{2}, \quad \forall a \geq a_0.$$

For any measurable  $A$ ,  $a \geq a_0$ ,

$$\begin{aligned} \sup_{f \in \mathcal{H}} \int 1_A |f| d\mu &\leq \sup_{f \in \mathcal{H}} \int_{\{|f| < a\}} 1_A |f| d\mu + \sup_{f \in \mathcal{H}} \int_{\{|f| \geq a\}} 1_A |f| d\mu \\ &\leq a\mu(A) + \sup_{f \in \mathcal{H}} \int_{\{|f| \geq a\}} |f| d\mu \leq a\mu(A) + \frac{\varepsilon}{2}. \end{aligned}$$

That  $\mathcal{H}$  is bounded in  $L_1$  follows by setting  $A = E$  and using the fact that  $\mu$  is finite. Fix  $a \geq a_0$ . For any  $A$  with  $\mu(A) \leq \varepsilon/(2a)$ , we get that  $\sup_{f \in \mathcal{H}} \int 1_A |f| d\mu$  is bounded from above by  $\varepsilon$ , hence the absolute continuity.

2. Suppose that  $\mathcal{H}$  is absolutely continuous and bounded in  $L_1$ . Denote the uniform  $L_1$  bound of  $\mathcal{H}$  by  $M$ . By Markov inequality,  $\forall a > 0$ ,

$$\mu(|f| > a) \leq \frac{1}{a} \int |f| d\mu \leq \frac{1}{a} M, \quad \forall f \in \mathcal{H}.$$

$\forall \varepsilon > 0$ , by absolute continuity,  $\sup_{f \in \mathcal{H}} \int 1_A |f| d\mu < \varepsilon$  as soon as  $\mu(A)$  is less than some  $\delta > 0$ . Fix  $a$  with  $M/a < \delta$ . Then setting  $A = \mu(|f| > a)$  gives the uniform integrability.  $\square$

**Thm 45** (Vitali convergence theorem). *Suppose that  $\mu$  is finite,  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ .*

(1) *If  $\{f_n\}$  is uniformly integrable, then  $f \in L_1$  and*

$$(6.10) \quad \int f_n d\mu \rightarrow \int f d\mu.$$

(2)  *$f_n, f$  are nonnegative integrable, then (6.10) implies that  $\{f_n\}$  is uniformly integrable.*

PROOF. The proof is given for  $f_n \xrightarrow{a.e.} f$ .

1. If  $f_n$  is uniformly integrable, then  $f$  is integrable by Theorem 44 and Fatou's lemma. Define

$$f_{n,a} = 1_{\{|f_n| < a\}} f_n, \quad f_a = 1_{\{|f| < a\}} f.$$

It follows that  $f_{n,a} \rightarrow f_a$ , *a.e.* provided  $\mu(|f| = a) = 0$ . By bounded dominated convergence,

$$\int f_{n,a} d\mu \rightarrow \int f_a d\mu.$$

Writing

$$(6.11) \quad \int_{\{|f_n| \geq a\}} f_n d\mu = \int f_n d\mu - \int f_{n,a} d\mu$$

and

$$(6.12) \quad \int_{\{|f| \geq a\}} f d\mu = \int f d\mu - \int f_a d\mu,$$



we see that

$$\begin{aligned}
& \limsup_n \left| \int f_n d\mu - \int f d\mu \right| \\
& \leq \limsup_n \left| \int f_{n,a} d\mu - \int f_a d\mu \right| + \sup_n \int_{\{|f_n| \geq a\}} |f_n| d\mu + \int_{\{|f| \geq a\}} |f| d\mu \\
& = \sup_n \int_{\{|f_n| \geq a\}} |f_n| d\mu + \int_{\{|f| \geq a\}} |f| d\mu.
\end{aligned}$$

Note  $\mu(|f| = a) = 0$  for all but countably many  $a$ . Sending  $a \rightarrow \infty$  proves (6.10).

**2.** Suppose  $f_n, f$  are nonnegative integrable and (6.10) holds. Write

$$\int_{\{|f_n| \geq a\}} f_n d\mu = \int_{\{|f| \geq a\}} f d\mu + \left( \int_{\{|f_n| \geq a\}} f_n d\mu - \int_{\{|f| \geq a\}} f d\mu \right).$$

Since  $f$  is integrable, the first term is less than  $\varepsilon/2$  when  $a$  is larger than some  $a_0$ . If  $\mu(|f| = a) = 0$ , (6.11) and (6.12) indicate the term in the bracket is also less than  $\varepsilon/2$  when  $n$  is larger than some  $n_0$ .

Therefore,

$$\sup_{n > n_0} \int_{\{|f_n| \geq a\}} f_n d\mu \leq \varepsilon, \quad \forall a > a_0 \text{ with } \mu(|f| = a) = 0.$$

Since the finite family  $\{f_1, \dots, f_{n_0}\}$  is uniformly integrable, the uniform integrability of  $\{f_n, n \geq 1\}$  follows.  $\square$

Additional details on the proof of Theorem 45. Suppose  $|f_n(x)| \rightarrow |f(x)| < a$ . Then for large  $n$ ,  $|f_n(x)| < a$ . So  $1_{\{|f_n| < a\}}$  and  $1_{\{|f| < a\}}$  are both equal to 1, it follows  $f_{n,a} \rightarrow f_a$  at  $x$ . The same is true for  $x$  with  $|f(x)| > a$ . If  $|f(x)| = a \neq 0$ , then  $f_{n,a}(x) \rightarrow f_a(x)$  may not happen, since in this case  $f_a(x) = 0$  while there could be a subsequence  $n_k$  with  $f_{n_k}(x) < a$  so that

$$f_{n_k,a}(x) = f_{n_k}(x) \rightarrow f(x) \neq 0.$$

But if the set  $\{x : |f(x)| = a\}$  has zero  $\mu$ -measure, then  $f_{n,a} \rightarrow f_a$ , *a.e.* Fortunately the set of  $a$  for which  $\mu(|f| = a)$  is not zero is at most countable. Indeed, let

$$F(x) = \mu(|f| \leq x).$$

Then  $F(x)$  is non-decreasing, hence has at most countably many discontinuities.  $F$  is (right-continuous and thus) discontinuous at  $x = a$  if and only if

$$\mu(|f| = a) = F(a) - F(a-) \neq 0.$$

This verifies that  $\mu(|f| = a) = 0$  for all but countably many  $a$ .

**COROLLARY 2.** *Suppose that  $\mu$  is finite,  $f_n, f$  are integrable. If  $f_n \xrightarrow{a.e.} f$  or  $f_n \xrightarrow{\mu} f$ , then these are equivalent:*

(1)  $\{f_n\}$  is uniformly integrable;

(2)  $\int |f_n - f| d\mu \rightarrow 0;$

(3)  $\int |f_n| d\mu \rightarrow \int |f| d\mu.$

## 6.5. Summary of various convergences.

$$f_n \xrightarrow{\mu} f \quad \begin{array}{c} \text{sub} \text{ subseq} \\ \Longleftrightarrow \\ \text{Thm 20} \end{array} \quad f_n \xrightarrow{a.u.} f$$

Markov  $\Uparrow$

$$f_n \xrightarrow{L_p} f \quad \begin{array}{c} \text{has a subseq} \\ \Longrightarrow \\ \Longleftarrow \\ \|f_n\|_p \rightarrow \|f\|_p \\ \text{Thm 40} \end{array} \quad f_n \xrightarrow{a.e.} f \quad \begin{array}{c} \xLeftrightarrow{\mu \text{ finite}} \\ \Longleftarrow \\ \text{Thm 19} \end{array} \quad f_n \xrightarrow{a.u.} f \quad \Downarrow \text{Thm 18}$$

$$f_n \xrightarrow{d} f \quad \begin{array}{c} \text{in a prob space} \\ \Longleftarrow \\ \text{Thm 24} \end{array} \quad f_n \xrightarrow{\mu} f$$

## 7. 概率空间的积分

**7.1. Expected value.**  $(\Omega, \mathcal{F}, P)$  is a probability space,  $X$  a r.v.

**Def 30.** *Expectation, written  $EX$ ,*

$$EX = \int X dP.$$

Suppose  $X$  is discrete, i.e.,  $X$  takes values in a finite or infinitely countable *distinct* sequence  $\{x_1, x_2, \dots\}$ . Then its expectation  $(\int X dP$  computed according to (5.1)) equals

$$EX = \sum_i x_i P(X = x_i).$$

The mapping  $i \mapsto P(X = x_i)$  is called the probability mass function of  $X$ . If  $Y = g(X)$  for some measurable function  $g$ , then  $Y$  is discrete with values in, say,  $\{y_1, y_2, \dots\}$ . The expectation of  $Y$ , computed in the

same way as  $EX$ , is

$$EY = \sum_i y_i P(Y = y_i).$$

To calculate  $EY$ , we first need to find its probability mass function  $i \mapsto P(Y = y_i)$ . This can be complicated, and it is avoided by using the "*law of the unconscious statistician*",

$$EY = \sum_i g(x_i) P(X = x_i).$$

This turns out to be a change of variables formula (see also Theorem [48](#)).

**Thm 46** (Change of variables formula). *Let  $(\Omega, \mathcal{F}, P)$  be a probability space,  $X$  a r.v, and  $g \in \mathcal{B}_{\mathbb{R}}$ . If  $g \geq 0$  or  $\int_{\Omega} |g(X)| dP < \infty$ , then*

$$(7.1) \quad Eg(X) = \int_{\Omega} g(X) dP = \int_{\mathbb{R}} g(x) d\mu_X.$$

Here  $\mu_X(A) = PX^{-1}(A) = P(X \in A)$ ,  $\forall A \in \mathcal{B}_{\mathbb{R}}$  is the probability induced by  $X$  (Section 5.3), which will be called the **distribution** of  $X$ .

**PROOF. 1.** The nonnegative case  $g \geq 0$ . If  $g = 1_A$ , then  $g(X(\omega)) = 1_A(X(\omega)) = 1_{X^{-1}(A)}(\omega)$ , so (7.1) reduces to the definition of  $\mu_X$ . By linearity, (7.1) holds for simple functions. If  $g_n$  are simple functions such that  $0 \leq g_n(x) \uparrow g(x)$ , then  $0 \leq g_n(X(\omega)) \uparrow g(X(\omega))$ , then (7.1) follows by monotone convergence theorem.

**2.** The case  $\int_{\Omega} |g(X)| dP < \infty$ . Applying step 1 to  $|g(X)|$  shows that  $g$  is integrable with respect to  $\mu_X$ , hence the integrability of  $g^+$ ,  $g^-$ , and (7.1) follows from subtracting  $Eg^-(X) = \int_{\mathbb{R}} g^-(x) d\mu_X$  from  $Eg^+(X) = \int_{\mathbb{R}} g^+(x) d\mu_X$ . □

The probability  $\mu_X$  equals (as a result of the uniqueness Theorem 10) the measure  $\mu$  constructed from the distribution function  $F$

of  $X : \mu((a, b]) = F(b) - F(a), \forall a, b$ . The measure  $\mu$  is called a Lebesgue-Stieltjes measure and its integral is the Lebesgue-Stieltjes integral (section 7.3). The above formula thus relates integral on a probability space to Lebesgue-Stieltjes integral over  $\mathbb{R}$ . The rightmost term of (7.1) is also written as  $\int g dF$ , i.e.

$$Eg(X) = \int_{\mathbb{R}} g(x) dF.$$

**REMARK 2.** *An implication of Theorem 46 is that the integration (e.g. the expectation and variance) of a random variable is a distributional property, i.e., it depends on the random variable only through its distribution. This lays the basis for applying probability theory tools such as Skorohod Theorem (Theorem 26).*

**Def 31.** *Variance, written  $\text{Var}(X)$ ,*

$$\text{Var}(X) = \int (X - EX)^2 dP = E(X - EX)^2.$$



It is easy to see that

$$\text{Var}(X) = EX^2 - (EX)^2.$$

**Def 32.** *k-th moment,  $k = 1, 2, \dots$ ,*

$$E(X^k) = \int X^k dP.$$

**Example 14** (Bernoulli distribution). *Let  $0 < p < 1$ .  $X \sim \text{Bernoulli}(p)$  if  $P(X = 1) = p$ ,  $P(X = 0) = 1 - p$ . Then*

$$EX = 1 \cdot p + 0 \cdot (1 - p) = p.$$

$$\text{Var}(X) = EX^2 - (EX)^2 = p - p^2 = p(1 - p).$$

**Example 15** (Poisson distribution). *Let  $\lambda > 0$ .  $X \sim \text{Poisson}(\lambda)$  if*

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

*Then*

$$EX = \sum_{k=0}^{\infty} k \cdot e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^k}{(k-1)!} = \lambda.$$

$$E(X(X-1)) = \sum_{k=0}^{\infty} k(k-1) \cdot e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \sum_{k=2}^{\infty} \frac{\lambda^k}{(k-2)!} = \lambda^2.$$

Hence  $EX^2 = \lambda^2 + \lambda$ , and

$$\text{Var}(X) = (\lambda^2 + \lambda) - \lambda^2 = \lambda.$$

**Example 16 (Geometric distribution).** *Repeatedly flip a coin with head probability  $p$  and stop only when the head appears. The number of tosses  $X$  has the distribution*

$$P(X = k) = (1-p)^{k-1}p, \quad k = 1, 2, \dots$$

*The distribution of  $X$  is called geometric, denoted by  $X \sim \text{Geom}(p)$ ,*

$$EX = \frac{1}{p}, \quad \text{Var}(X) = \frac{1-p}{p^2}.$$

**7.2. Properties of expectation.**  $X, Y$  are random variables. The following are immediate from section 6.1.

**Jensen inequality:** if  $X$  integrable,  $\varphi$  convex, then

$$\varphi(EX) \leq E\varphi(X).$$

**Hölder inequality:** if  $p, q \geq 1, 1/p + 1/q = 1$ , then

$$E|XY| \leq \|X\|_p \|Y\|_q.$$

**Minkowski inequality:** if  $p \geq 1$ , then

$$\|X + Y\|_p \leq \|X\|_p + \|Y\|_p.$$

**Thm 47.**  $0 < s < t < \infty$ ,  $X$  is a r.v. Then  $\|X\|_s \leq \|X\|_t$ .

PROOF. By Hölder inequality with  $p = \frac{t}{s}$ ,  $q = \frac{t}{t-s}$ ,

$$\|X\|_s^s = E|X|^s \leq (E|X|^{sp})^{1/p} (E1^q)^{1/q} = (E|X|^t)^{s/t} = \|X\|_t^s.$$

□

**Example 17.** If  $X$  has  $EX^2 < \infty$ , then its expectation and variance exist, since  $E|X| \leq \|X\|_2 < \infty$ , and

$$0 \leq \text{Var}(X) \leq EX^2.$$

**7.3. Lebesgue-Stieltjes and Riemann-Stieltjes integrals.** Let  $G$  be a **generalized distribution function**, i.e., nondecreasing, right-continuous on  $\mathbb{R}$ . There is a unique measure  $\mu$  such that

$$(7.2) \quad \mu((a, b]) = G(b) - G(a), \quad \forall a, b.$$

The measure  $\mu$  constructed this way is called a **Lebesgue-Stieltjes measure**. Integration with respect to Lebesgue-Stieltjes measure is called **Lebesgue-Stieltjes integral**, denoted by  $\int f d\mu$  or  $\int f dG$ .

REMARK 3. *Under suitable conditions (see below),  $\int f dG$  may be interpreted as Riemann-Stieltjes integral. Since this does not provide anything new in the context of general measure theory,  $\int f dG$  is best understood as a notional variant of  $\int f d\mu$ , and hence by convention (see (7.2))  $\int_a^b f dG$  means  $\int_{(a,b]} f dG$  or  $\int_{(a,b]} f d\mu$ .*

Here we recall a few facts about Riemann-Stieltjes integration. Let  $G$  be the function as in (7.2),  $f$  a bounded function on  $[a, b]$ . Corresponding to each partition  $\mathcal{P} : a = x_0 < x_1 < \cdots < x_n = b$ , we consider

$$L(\mathcal{P}, f) = \sum_{i=1}^n \inf_{x \in [x_{i-1}, x_i]} f(x) \Delta G_i, \quad U(\mathcal{P}, f) = \sum_{i=1}^n \sup_{x \in [x_{i-1}, x_i]} f(x) \Delta G_i.$$

Here  $\Delta G_i = G(x_i) - G(x_{i-1})$ . Define

$$R_* f = \sup_{\mathcal{P}} L(\mathcal{P}, f), \quad R^* f = \inf_{\mathcal{P}} U(\mathcal{P}, f).$$

If  $R_* f = R^* f$ , then  $f$  is Riemann-Stieltjes integrable with respect to  $G$ , the common value, written  $(R-S) \int f$ , is called the Riemann-Stieltjes integral. For simplicity we have omitted the dependence of the integral on  $G$  in the notations.

A sufficient condition for Riemann-Stieltjes integrability is this: Suppose  $f$  is bounded on  $[a, b]$ , has at most finitely many discontinuities,  $G$  is continuous at every point where  $f$  is discontinuous. Then  $f$  is Riemann-Stieltjes integrable with respect to  $G$ .

**Example 18.** *If  $a < s < b$ ,  $f$  is bounded on  $[a, b]$ , continuous at  $s$  and  $G(x) = 1_{[s, \infty)}(x)$ . Then*

$$(R-S) \int_a^b f dG = f(s).$$

*Indeed, consider partitions  $\mathcal{P} = \{x_0, x_1, x_2, x_3\}$ ,  $a = x_0$  and  $x_1 < x_2 = s < x_3 = b$ . Then  $\Delta G_2 = 1$ ,  $\Delta G_i = 0$  if  $i \neq 2$ ,*

$$L(\mathcal{P}, f) = \inf_{x \in [x_1, x_2]} f(x), \quad U(\mathcal{P}, f) = \sup_{x \in [x_1, x_2]} f(x).$$

*Since  $f$  is continuous at  $s$ , we see that  $L(\mathcal{P}, f)$  and  $U(\mathcal{P}, f)$  converge to  $f(s)$  as  $x_1 \rightarrow s$ .*

**Thm 48.** Suppose  $c_n \geq 0$ ,  $\sum c_n < \infty$ ,  $\{s_n\}$  is a sequence of distinct points in  $(a, b)$ , and

$$G(x) = \sum_{n=1}^{\infty} c_n 1_{[s_n, \infty)}(x).$$

If  $f$  is continuous on  $[a, b]$ , then

$$(R-S) \int_a^b f dG = \sum_{n=1}^{\infty} c_n f(s_n).$$

PROOF. Exercise. □

If we denote by  $L_* f$  the integral in (5.1) with the  $G$ -induced Lebesgue-Stieltjes measure in the role of  $\mu$ , and by  $L^* f$  the integral in (5.2). Then

$$R_* f \leq L_* f \leq L^* f \leq R^* f.$$

Therefore if, for instance,  $f$  is continuous on  $[a, b]$ , then it is Riemann-Stieltjes integrable, hence Lebesgue-Stieltjes integrable.

## 7.4. $L_p$ convergence and uniform integrability.

**Thm 49.**  $(\Omega, \mathcal{F}, P)$  is a probability space,  $0 < p < \infty$ ,  $X_n \in L_p$ ,  $X \in \mathcal{F}$ . If  $X_n \xrightarrow{P} X$ , then these are equivalent:

- (1)  $\{|X_n|^p\}$  is uniformly integrable;
- (2)  $X \in L_p$ ,  $E(|X_n - X|^p) \rightarrow 0$ ;
- (3)  $X \in L_p$ ,  $E(|X_n|^p) \rightarrow E(|X|^p)$ .

**PROOF. 1.** Suppose that  $\{|X_n|^p\}$  is uniformly integrable. Observe that  $X \in L_p$  by Theorem 45, hence  $\{|X_n - X|^p\}$  is uniformly integrable since  $|X_n - X|^p \leq C_p(|X_n|^p + |X|^p)$  where  $C_p = 2^{p-1} \vee 1$ . Note also that  $|X_n - X|^p \xrightarrow{P} 0$ . Therefore (1) implies (2) is a consequence of Theorem 45 with  $f_n = |X_n - X|^p$ .

**2.** (2) implies (3) because  $\left| \|X_n\|_p - \|X\|_p \right| \leq \|X_n - X\|_p$ ,  $0 < p < \infty$  (Theorem 36, Theorem 38).

**3.** (3) implies (2) follows from an application of Theorem 45 with  $f_n = |X_n|^p$ . □



We notice another criterion for uniform integrability, in addition to Theorem 44.

LEMMA 50. *Let  $(\Omega, \mathcal{F}, P)$  be a probability space,*

$$\mathcal{H} = \{X_t : t \in T, E|X_t| < \infty\}.$$

*Suppose that  $g \geq 0$  is an increasing function on  $[0, \infty)$  such that*

$$\lim_{s \rightarrow \infty} \frac{g(s)}{s} = \infty$$

*and*

$$\sup_{X \in \mathcal{H}} \int g(|X|) dP < \infty.$$

*Then  $\mathcal{H}$  is uniformly integrable.*

PROOF.  $\forall \varepsilon > 0$ . Fix  $a > 0$  so that

$$\frac{1}{a} \sup_{X \in \mathcal{H}} \int g(|X|) dP < \varepsilon.$$

There is  $s_0 > 0$  such that  $g(s) \geq as$  for all  $s \geq s_0$ . Hence,  $\forall X \in \mathcal{H}$ ,  $s \geq s_0$ ,

$$\int_{\{|X| \geq s\}} |X| dP \leq \frac{1}{a} \int_{\{|X| \geq s\}} g(|X|) dP \leq \frac{1}{a} \sup_{X \in \mathcal{H}} \int g(|X|) dP < \varepsilon.$$

□

## 8. 乘积测度空间

Let  $(X, \mathcal{X}, \mu)$ ,  $(Y, \mathcal{Y}, \nu)$  be measure spaces. The problem is to construct a **product measure**  $\pi$  on  $X \times Y$  such that

$$(8.1) \quad \pi(A \times B) = \mu(A)\nu(B) \text{ for } A \in \mathcal{X}, B \in \mathcal{Y}.$$

### 8.1. Product $\sigma$ -field.

**Def 33.** *In the product space  $X \times Y$ , a **measurable rectangle** is a product of the form*

$$A \times B, A \in \mathcal{X}, B \in \mathcal{Y}.$$

*Let*

$$(8.2) \quad \mathcal{S} = \{A \times B : A \in \mathcal{X}, B \in \mathcal{Y}\}$$

*be the class of measurable rectangles on  $X \times Y$ . The **product  $\sigma$ -field** on  $X \times Y$  is then defined as*

$$\mathcal{X} \times \mathcal{Y} = \sigma(\mathcal{S}).$$

The space  $X \times Y$  equipped with this product  $\sigma$ -field is called **product measurable space**.

As the example below shows, the product  $\sigma$ -field is generally larger than the class of measurable rectangles.

**Example 19.**  $\mathcal{B}(\mathbb{R}^2) = \mathcal{B}(\mathbb{R}) \times \mathcal{B}(\mathbb{R})$ . Here  $\mathcal{B}(\mathbb{R}^2)$  is the usual  $\sigma$ -field on  $\mathbb{R}^2$  generated by the class of products of one-dimensional intervals.

**Example 20.** If  $A \times B \in \mathcal{S}$ , then

$$(A \times B)^c = A^c \times Y + A \times B^c \in \mathcal{S}.$$

From this it is easy to check that  $\mathcal{S}$  is a semi-ring and  $X \times Y \in \mathcal{S}$ .

**Def 34.** The **section of a set**  $E \in X \times Y$  at  $x \in X$  is

$$E_x = \{y : (x, y) \in E\}.$$

Similarly  $E_y = \{x : (x, y) \in E\}$  is the section at  $y \in Y$ . The **section of a function**  $f(x, y)$  at  $x \in X$  is the mapping

$$y \mapsto f(x, y).$$

The section at  $y \in Y$  is  $x \mapsto f(x, y)$ .

**Example 21.** If  $E, E_k \in \mathcal{X} \times \mathcal{Y}$ ,  $x \in X$ , then

$$(E^c)_x = (E_x)^c, \quad \left(\bigcup_k E_k\right)_x = \bigcup_k (E_k)_x, \quad \left(\bigcap_k E_k\right)_x = \bigcap_k (E_k)_x$$

**Thm 51.** (1) Sections of  $\mathcal{X} \times \mathcal{Y}$ -measurable set are measurable.

(2) Sections of  $\mathcal{X} \times \mathcal{Y}$ -measurable function are measurable.

PROOF. 1. Fix  $x \in X$ . Consider the mapping  $R_x : Y \mapsto X \times Y$  defined by  $R_x(y) = (x, y)$ . We intend to prove that  $R_x$  is  $\mathcal{Y}$ -measurable so that the conclusion follows immediately:

$$E_x = R_x^{-1}E \in \mathcal{Y}, \quad \forall E \in \mathcal{X} \times \mathcal{Y}.$$

If  $E = A \times B$  is a measurable rectangle, then  $R_x^{-1}E = B \in \mathcal{Y}$ . This shows that  $R_x$  is  $\mathcal{Y}$ -measurable, since  $\mathcal{X} \times \mathcal{Y}$  is generated by measurable rectangles. So the first part is proved.

2. If  $f \in \mathcal{X} \times \mathcal{Y}$ , then  $f(x, \cdot) = f \circ R_x(\cdot)$  is  $\mathcal{Y}$ -measurable by measurable composition.

3. The conclusion for fixed  $y \in Y$  is proved similarly. □

**8.2. Product measure space.** Let  $(X, \mathcal{X}, \mu), (Y, \mathcal{Y}, \nu)$  be measure spaces.

LEMMA 52. *Suppose that  $\mu$  and  $\nu$  are finite. If  $E \in \mathcal{X} \times \mathcal{Y}$ , then the mapping  $x \mapsto \nu(E_x)$  is  $\mathcal{X}$ -measurable,  $y \mapsto \mu(E_y)$  is  $\mathcal{Y}$ -measurable.*

PROOF. Let  $\mathcal{L}$  be the class of  $E \in \mathcal{X} \times \mathcal{Y}$  that has the stated property. Then  $\mathcal{L}$  is a  $\lambda$ -system. Indeed, it is easy to see that  $X \times Y \in \mathcal{L}$ . If  $E, F \in \mathcal{L}$  with  $E \subset F$ , then

$$x \mapsto \nu((F \setminus E)_x) = \nu(F_x \setminus E_x) = \nu(F_x) - \nu(E_x)$$

is  $\mathcal{X}$ -measurable (the finiteness of  $\nu$  is used to justify subtraction). If  $E_k \in \mathcal{L}$ ,  $E_k \subset E_{k+1}$ , then

$$x \mapsto \nu\left(\left(\bigcup_k E_k\right)_x\right) = \nu\left(\bigcup_k (E_k)_x\right) = \lim_k \nu((E_k)_x)$$

is  $\mathcal{X}$ -measurable. This shows that  $\mathcal{L}$  is a  $\lambda$ -system. For any measurable rectangle  $E = A \times B$ , the function

$$x \mapsto \nu(E_x) = 1_A(x)\nu(B)$$

is  $\mathcal{X}$ -measurable. So  $\mathcal{L}$  contains the  $\pi$ -system of measurable rectangles, thus coincides with  $\mathcal{X} \times \mathcal{Y}$  by  $\pi$ - $\lambda$  Theorem.  $\square$

LEMMA 53. *Let  $(X, \mathcal{X}, \mu)$ ,  $(Y, \mathcal{Y}, \nu)$  be  $\sigma$ -finite measure spaces. Define*

$$\pi_{21}(E) = \int_X \nu(E_x) d\mu, \quad E \in \mathcal{X} \times \mathcal{Y}$$

*and*

$$\pi_{12}(E) = \int_Y \mu(E_y) d\nu, \quad E \in \mathcal{X} \times \mathcal{Y}.$$

*Then  $\pi_{21}(E) = \pi_{12}(E)$  for  $E \in \mathcal{X} \times \mathcal{Y}$ . Moreover,  $\pi_{21}$ ,  $\pi_{12}$  satisfy (8.1).*

PROOF. 1. First suppose  $\mu, \nu$  are finite. Let  $\mathcal{L}$  be the class of  $E \in \mathcal{X} \times \mathcal{Y}$  such that  $\pi_{21}(E) = \pi_{12}(E)$ .  $\mathcal{L}$  contains all measurable

rectangle  $A \times B$ , since

$$\begin{aligned}\pi_{21}(A \times B) &= \int_X 1_A(x) \nu(B) d\mu = \mu(A) \nu(B) \\ &= \int_Y \mu(A) 1_B(y) d\nu = \pi_{12}(A \times B).\end{aligned}$$

It is easy to check that  $\mathcal{L}$  is a  $\lambda$ -system, hence equals  $\mathcal{X} \times \mathcal{Y}$  by  $\pi$ - $\lambda$  Theorem.

**2.** Now suppose  $\mu, \nu$  are  $\sigma$ -finite, then there are  $\{A_m\}, \{B_n\}$  that partition  $X$  and  $Y$  into disjoint sets of finite measure. Define

$$\mu_m(E) = \mu(E \cap A_m), \quad \nu_n(E) = \nu(E \cap B_n), \quad E \in \mathcal{X} \times \mathcal{Y}.$$

Step 1 is valid for these finite measures,

$$(8.3) \quad \pi_{21}^{(mn)}(E) = \pi_{12}^{(mn)}(E), \quad E \in \mathcal{X} \times \mathcal{Y},$$

where

$$\pi_{21}^{(mn)}(E) = \int_X \nu_n(E_x) d\mu_m,$$



and

$$\pi_{12}^{(mn)}(E) = \int_Y \mu_m(E_y) d\nu_n.$$

In addition, for measurable rectangle  $A \times B$ ,

$$(8.4) \quad \pi_{21}^{(mn)}(A \times B) = \mu_m(A) \nu_n(B) = \pi_{12}^{(mn)}(A \times B).$$

From Lemma 52,  $x \mapsto \nu_n(E_x)$  is measurable. Since  $\nu = \sum_n \nu_n$ ,  $x \mapsto \nu(E_x)$  is measurable. The same can be said for  $y \mapsto \mu(E_y)$ . Therefore  $\pi_{21}$ ,  $\pi_{12}$  are well-defined for the  $\sigma$ -finite case. By Example 10, 11 and (8.3), for  $E \in \mathcal{X} \times \mathcal{Y}$ ,

$$\begin{aligned} \pi_{21}(E) &= \int_X \nu(E_x) d\mu = \sum_m \int_X \nu(E_x) d\mu_m = \sum_m \int_X \sum_n \nu_n(E_x) d\mu_m \\ &= \sum_{m,n} \pi_{21}^{(mn)}(E) = \sum_{m,n} \pi_{12}^{(mn)}(E) = \sum_n \int_Y \sum_m \mu_m(E_y) d\nu_n \\ &= \pi_{12}(E). \end{aligned}$$

Particularly, this together with (8.4) yields that, for measurable rectangle  $A \times B$ ,

$$\pi_{21}(A \times B) = \sum_{m,n} \mu_m(A) \nu_n(B) = \pi_{12}(A \times B).$$

□

**Thm 54.** *Let  $(X, \mathcal{X}, \mu)$ ,  $(Y, \mathcal{Y}, \nu)$  be  $\sigma$ -finite measure spaces. Then*

$$\pi(E) \triangleq \pi_{21}(E) = \pi_{12}(E)$$

*defines the unique  $\sigma$ -finite measure on  $X \times Y$  that satisfies (8.1).*

**PROOF.** Decompose  $X$  and  $Y$  into disjoint sets  $\{A_m\}$ ,  $\{B_n\}$  of finite measure, and define  $\mu_m$ ,  $\nu_n$  as before. Then  $X \times Y$  is the disjoint union of  $\{A_m \times B_n\}$  and each  $A_m \times B_n$  has finite  $\pi$ -measure:  $\pi(A_m \times B_n) = \mu_m(A_m) \nu_n(B_n)$ . It follows that  $\pi$  is  $\sigma$ -finite. The uniqueness is a consequence of Theorem 10, since measurable rectangles form a  $\pi$ -system (Example 20). □

In the future, the product measure  $\pi$  of  $\mu$  and  $\nu$  will be written as  $\mu \times \nu$ .

### 8.3. Fubini's Theorem.

**Thm 55.** *Let  $(X, \mathcal{X}, \mu)$ ,  $(Y, \mathcal{Y}, \nu)$  be  $\sigma$ -finite measure spaces,  $\pi$  the product measure constructed in Theorem 54,  $f$  a  $\mathcal{X} \times \mathcal{Y}$ -measurable function. If  $f$  is nonnegative or  $\int_{X \times Y} |f| d\pi < \infty$ . Then*

$$\int_X \left[ \int_Y f(x, y) d\nu \right] d\mu = \int_Y \left[ \int_X f(x, y) d\mu \right] d\nu = \int_{X \times Y} f d\pi.$$

*It is implicit in the statement that all integrands are integrable. In the nonnegative case, if one of the above integrals is infinite, so it is with the other two.*

**PROOF.** The conclusion holds for measurable indicator function by Theorem 54, and hence simple function by linearity of integration.

Then the monotone convergence theorem gives the conclusion for non-negative measurable function. If  $\int_{X \times Y} |f| d\pi < \infty$ , then applying the nonnegative case to  $|f|$ ,

$$\int_X \left[ \int_Y |f| d\nu \right] d\mu = \int_{X \times Y} |f| d\pi < \infty.$$

It follows that

$$\int_Y |f| d\nu < \infty, \text{ a.e. } x.$$

Hence it makes sense (outside a set of zero  $\mu$ -measure) to write

$$\int_Y f d\nu = \int_Y f^+ d\nu - \int_Y f^- d\nu.$$

Now the desired property follows by integrating over  $X$  and using the result for nonnegative integrand. The same reasoning applies to

$$\int_Y \left[ \int_X |f| d\mu \right] d\nu.$$

□

## 8.4. Applications.

**Example 22 (Euler-Poisson integral).** Let  $I = \int_{\mathbb{R}} e^{-x^2} dx$ . By Fubini's theorem

$$I^2 = \iint_{\mathbb{R}^2} e^{-(x+y)^2} dx dy = \iint_{\substack{r \geq 0, \\ 0 \leq \theta < 2\pi}} e^{-r^2} r dr d\theta.$$

Again by Fubini's theorem, the double integral on the RHS is written as an iterated integral and evaluated to give  $I^2 = \pi$ , so

$$I = \int_{\mathbb{R}} e^{-x^2} dx = \sqrt{\pi}.$$

**Thm 56.** Let  $(X, \mathcal{X}, \mu)$  be a  $\sigma$ -finite measure space,  $f \geq 0$  measurable. Then

$$\int f d\mu = \int_0^\infty \mu(f \geq t) dt = \int_0^\infty \mu(f > t) dt$$

PROOF. Since  $f$  is nonnegative, we may write (recall our convention Remark 3),  $\forall x$ ,

$$f(x) = \int_0^{f(x)} dt = \int_0^\infty 1_{(0, f(x)]}(t) dt.$$

Notice that

$$1_{(0, f(x)]}(t) = 1_{\{x: f(x) \geq t\}}(x).$$

Then using Fubini theorem

$$\begin{aligned} \int f d\mu &= \int \int_0^\infty 1_{(0, f(x)]}(t) dt d\mu = \int_0^\infty \int 1_{(0, f(x)]}(t) d\mu dt \\ &= \int_0^\infty \int 1_{\{x: f(x) \geq t\}}(x) d\mu dt = \int_0^\infty \mu(\{x : f(x) \geq t\}) dt. \end{aligned}$$

Since the set of  $t$  such that  $\mu(\{x : f(x) = t\})$  is non-zero is at most countable, hence has zero Lebesgue measure. Thus the two integrals are equal,

$$\int_0^\infty \mu(\{x : f(x) \geq t\}) dt = \int_0^\infty \mu(\{x : f(x) > t\}) dt.$$

□

If  $f$  takes values in  $\{y_1, y_2, \dots\}$  and  $0 \leq y_1 < y_2 < \dots$ . Then  $t \mapsto \mu(f \geq t)$  is a step function

$$\mu(f \geq t) = \begin{cases} \mu(f \geq y_1), & 0 \leq t \leq y_1; \\ \mu(f \geq y_n), & y_{n-1} < t \leq y_n. \end{cases}$$

Hence Theorem 56 reduces to

$$\int f d\mu = y_1 \mu(f \geq y_1) + \sum_{n=2}^{\infty} (y_n - y_{n-1}) \mu(f \geq y_n).$$

A particular case of this is  $f$  taking values in nonnegative integers.

**COROLLARY 3.** *Let  $(X, \mathcal{X}, \mu)$  be a  $\sigma$ -finite measure space,  $f$  measurable with values in  $\{0, 1, 2, \dots\}$ . Then*

$$\int f d\mu = \sum_{n=1}^{\infty} \mu(f \geq n) = \sum_{n=0}^{\infty} \mu(f > n).$$

**Example 23.** *Example 11 has validated interchanging the order of summation as an application of dominated convergence theorem. The same can also be proved directly using Fubini theorem.*

**Thm 57 (Integration by parts).** *Let  $F, G$  be two nondecreasing, right-continuous functions on  $\mathbb{R}$ , then, for  $a < b$ ,*

$$F(b)G(b) - F(a)G(a) = \int_{(a,b]} G(x)dF(x) + \int_{(a,b]} F(x-)dG(x),$$

*or equivalently*

$$\begin{aligned} F(b)G(b) - F(a)G(a) &= \int_{(a,b]} G(-x)dF(x) + \int_{(a,b]} F(x-)dG(x) \\ &\quad + \sum_{a < x \leq b} \Delta F(x)\Delta G(x), \end{aligned}$$

*where  $\Delta F(x) = F(x) - F(x-)$ .*

**PROOF.** Denote respectively by  $\mu, \nu$  the Lebesgue-Stieltjes measure induced by  $F, G$ . Let  $\pi = \mu \times \nu$  be the product measure of  $\mu, \nu$ .



Using Fubini Theorem we have

$$\begin{aligned}
 \pi((x, y) : a < x < y \leq b) &= \int_{(a,b]} \int_{(a,y)} d\mu(x) d\nu(y) \\
 &= \int_{(a,b]} [F(y-) - F(a)] d\nu(y) \\
 &= \int_{(a,b]} F(y-) d\nu(y) - F(a)[G(b) - G(a)]
 \end{aligned}$$

and similarly

$$\begin{aligned}
 \pi((x, y) : a < y \leq x \leq b) &= \int_{(a,b]} \int_{(a,x]} d\nu(y) d\mu(x) \\
 &= \int_{(a,b]} [G(x) - G(a)] d\mu(x) \\
 &= \int_{(a,b]} G(x) d\mu(x) - G(a)[F(b) - F(a)]
 \end{aligned}$$

By the construction of  $\pi$ ,

$$(F(b) - F(a))(G(b) - G(a)) = \pi((x, y) \in (a, b] \times (a, b]).$$

The first conclusion follows by putting together these equations. To complete the proof, it suffices to note that

$$\begin{aligned} \pi((x, y) : a < y = x \leq b) &= \int_{(a, b]} \nu(\{x\}) d\mu(x) \\ &= \int_{(a, b]} [G(x) - G(x-)] d\mu(x) \\ &= \sum_{a < x \leq b} \Delta F(x) \Delta G(x). \end{aligned}$$

□

**8.5. Finite-dimensional product space.** The discussion for two-dimensional product space obviously extends to finite dimensions. Denote by  $\mathcal{B}(\mathbb{R}^n)$  the  $\sigma$ -field on  $\mathbb{R}^n$  generated by open sets. The product

$\sigma$ -field  $\prod_{i=1}^n \mathcal{B}(\mathbb{R})$  is defined as being generated by measurable rectangles with sides in  $\mathcal{B}(\mathbb{R})$ . Similar to Example [19](#), we have on  $\mathbb{R}^n$ ,

$$(8.5) \quad \mathcal{B}(\mathbb{R}^n) = \prod_{i=1}^n \mathcal{B}(\mathbb{R}).$$

## 9. 独立性 Independence

Let  $(\Omega, \mathcal{F}, P)$  be a probability space. A subset in  $\mathcal{F}$  is called an **event**, and an element of  $\Omega$  is called a **sample**.

### 9.1. Independence of events and random variables.

**Def 35.** *The events  $A$  and  $B$  are independent if*

$$P(A \cap B) = P(A)P(B).$$

**Thm 58.**  *$A, B$  are independent if and only if one of the three pairs are independent: (i)  $A^c, B$ ; (ii)  $A, B^c$ ; (iii)  $A^c, B^c$ .*

**Example 24.** *An event  $A$  is independent of any event if and only if  $P(A) = 0$  or  $1$ . In particular,  $\Omega$  and  $\emptyset$  are independent of any event.*

**Example 25 (Gambler's fallacy).** *In some situations, an individual erroneously think that certain event is more or less likely to happen in the future based on the outcome of the past events. This incorrect belief may lead a gambler in a coin flipping game to believe that after 100 successive heads, the next toss would be more likely to*

come up tail. The fallacy roots from the ignorance of the independence between tosses.

**Def 36.**  $\sigma$ -fields  $\mathcal{F}, \mathcal{G}$  are independent if

$$P(A \cap B) = P(A)P(B), \quad \forall A \in \mathcal{F}, \quad B \in \mathcal{G}.$$

**Def 37.** Random variables  $X, Y$  are independent if  $\sigma(X)$  and  $\sigma(Y)$  are independent, where  $\sigma(X) = X^{-1}(\mathcal{B}(\mathbb{R}))$ .

**Thm 59.**  $A, B$  are independent if and only if  $1_A$  and  $1_B$  are independent.

PROOF. Note that  $\sigma(1_A) = \{A, A^c, \emptyset, \Omega\}$  and the same for  $\sigma(1_B)$ . □

**Def 38.** A family of events  $\{A_1, \dots, A_n\}$  is independent if for any  $I \subset \{1, \dots, n\}$ ,

$$P\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} P(A_i).$$

Pairwise independence is weaker than independence.

**Example 26.** *Flip a fair coin twice and consider the events,*

$$A_1 = \{ \text{head-head, head-tail} \},$$

$$A_2 = \{ \text{head-head, tail-head} \},$$

$$A_3 = \{ \text{head-head, tail-tail} \}.$$

*Then  $A_1, A_2, A_3$  are pairwise independent but not independent, since*

$$P(A_i \cap A_j) = P(\text{head-head}) = \frac{1}{4} = P(A_i)P(A_j), \quad i \neq j.$$

$$P(A_1 \cap A_2 \cap A_3) = P(\text{head-head}) = \frac{1}{4} \neq \frac{1}{8} = P(A_1)P(A_2)P(A_3).$$

**Def 39.**  $\mathcal{A}_1, \dots, \mathcal{A}_n$  are classes of sets. They are independent if for any  $I \subset \{1, \dots, n\}$ ,

$$P\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} P(A_i), \quad \forall A_i \in \mathcal{A}_i.$$

If we denote by  $\mathcal{A}_i'$  the class formed by augmenting  $\mathcal{A}_i$  with  $\Omega$ . Then it is easy to see that  $\mathcal{A}_1, \dots, \mathcal{A}_n$  are independent if and only if

$\mathcal{A}'_1, \dots, \mathcal{A}'_n$  are independent. Definition 39 is thus equivalent to the *full-product form*,

$$(9.1) \quad P\left(\bigcap_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i), \quad \forall A_i \in \mathcal{A}'_i.$$

This form may bring added convenience when independence is to be verified. Since  $\Omega$  is contained in  $\sigma$ -field, the independence of random variables can be defined in this full-product form.

**Def 40.**  $X_1, \dots, X_n$  are independent if  $\sigma(X_1), \dots, \sigma(X_n)$  are independent, i.e.,

$$P\left(\bigcap_{i=1}^n \{X_i \in A_i\}\right) = \prod_{i=1}^n P(X_i \in A_i), \quad \forall A_i \in \mathcal{B}(\mathbb{R}), \quad i = 1, \dots, n.$$

**Thm 60.** Suppose that  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$  are independent  $\pi$ -systems. Then  $\sigma(\mathcal{A}_1), \sigma(\mathcal{A}_2), \dots, \sigma(\mathcal{A}_n)$  are independent.

PROOF. 1. Clearly it suffices to show that  $\sigma(\mathcal{A}_1), \mathcal{A}_2, \dots, \mathcal{A}_n$  are independent, since the conclusion applies to itself and would yield that  $\sigma(\mathcal{A}_1), \sigma(\mathcal{A}_2), \dots, \mathcal{A}_n$  are independent, and so on.

2. Now we show that  $\sigma(\mathcal{A}_1), \mathcal{A}_2, \dots, \mathcal{A}_n$  are independent. Fix  $A_i \in \mathcal{A}_i, i = 2, \dots, n$ . Let  $E = \bigcap_{i=2}^n A_i$  and

$$\mathcal{L}_E = \{A \in \sigma(\mathcal{A}_1) : P(A \cap E) = P(A)P(E)\}.$$

Then  $\mathcal{A}_1 \subset \mathcal{L}_E$ . In view of Example 24,  $\Omega \in \mathcal{L}_E$ . If  $B_1, B_2 \in \mathcal{L}_E$  and  $B_1 \subset B_2$ , then

$$\begin{aligned} P((B_2 - B_1) \cap E) &= P(B_2 \cap E) - P(B_1 \cap E) \\ &= P(B_2)P(E) - P(B_1)P(E) \\ &= P(B_2 - B_1)P(E). \end{aligned}$$



Hence  $B_2 - B_1 \in \mathcal{L}_E$ . Finally let  $B_k \in \mathcal{L}_E$ ,  $B_k \subset B_{k+1}$ , then

$$\begin{aligned} P\left(\left(\bigcup_{k=1}^{\infty} B_k\right) \cap E\right) &= \lim_k P(B_k \cap E) \\ &= \lim_k P(B_k)P(E) = P\left(\bigcup_{k=1}^{\infty} B_k\right)P(E). \end{aligned}$$

Thus  $\bigcup_{k=1}^{\infty} B_k \in \mathcal{L}_E$ . Therefore  $\mathcal{L}_E$  is a  $\lambda$ -system and  $\sigma(\mathcal{A}_1) \subset \mathcal{L}_E$ . The desired conclusion follows.  $\square$

Since  $\mathcal{B}(\mathbb{R})$  is generated by the class  $\mathcal{S} = \{(-\infty, a] : a \in \mathbb{R}\}$ , the  $\sigma$ -field  $\sigma(X)$  generated by  $X$  equals  $\sigma(\{X \leq a : a \in \mathbb{R}\})$ , which together with Theorem 60 gives the following criterion for independence in terms of distribution functions.

**Thm 61.**  $X_1, \dots, X_n$  are independent if and only if  $\forall x_1, \dots, x_n \in \mathbb{R}$ ,

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = \prod_{i=1}^n P(X_i \leq x_i).$$

Note the class  $\{X_i \leq x_i : x_i \in \mathbb{R}\}$  may not contain but can approximate  $\Omega$ , so it is still legal to use the full-product form (9.1).

Recall from Theorem 46 that each random variable  $X$  induces a probability  $\mu$  on  $\mathbb{R}$ , which is called the distribution of  $X$ . When random vector  $(X_1, \dots, X_n)$  is involved, the same can be said. From Section 5.3, we see that the random vector induces a probability on  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ , where  $\mathcal{B}(\mathbb{R}^n)$  equals the  $\sigma$ -field generated by measurable rectangles with sides in  $\mathcal{B}(\mathbb{R})$  (see (8.5)). The induced probability on  $\mathcal{B}(\mathbb{R}^n)$ , denoted by  $P_{X_1, \dots, X_n}$ , satisfies,  $\forall B_1, \dots, B_n \in \mathcal{B}(\mathbb{R})$ ,

$$P_{X_1, \dots, X_n}(B_1 \times \cdots \times B_n) = P(X_1 \in B_1, \dots, X_n \in B_n).$$

$P_{X_1, \dots, X_n}$  is called the **joint distribution** of  $(X_1, \dots, X_n)$ . If  $B_i$  takes the form  $(-\infty, x_i]$ , then we get an associated mapping

$$F_{X_1, \dots, X_n} : (x_1, \dots, x_n) \mapsto P(X_1 \leq x_1, \dots, X_n \leq x_n)$$

which is called the **joint distribution function**<sup>2</sup> of the random vectors  $(X_1, \dots, X_n)$ . Moreover, the general change of variables formula from Section 5.3 tells us that, for measurable  $g : \mathbb{R}^n \mapsto \mathbb{R}$ ,

$$(9.2) \quad Eg(X_1, \dots, X_n) = \int_{\mathbb{R}^n} g(x_1, \dots, x_n) dP_{X_1, \dots, X_n}.$$

whenever one of the integrals exists.

**Example 27.** *The discrete random variables  $X \in \{x_1, x_2, \dots\}$  and  $Y \in \{y_1, y_2, \dots\}$  are independent if and only if*

$$(9.3) \quad P(X = x_i, Y = y_j) = P(X = x_i)P(Y = y_j), \quad \forall i, j.$$

---

<sup>2</sup>In fact, by uniqueness  $P_{X_1, \dots, X_n}$  equals the Lebesgue-Stieltjes measure determined by the joint distribution function on  $\mathcal{B}(\mathbb{R}^n)$ .

PROOF. 1. Suppose that (9.3) holds.  $\forall x, y \in \mathbb{R}$ ,

$$\begin{aligned} P(X \leq x, Y \leq y) &= \sum_{\substack{i: x_i \leq x \\ j: y_j \leq y}} P(X = x_i, Y = y_j) \\ &= \sum_{\substack{i: x_i \leq x \\ j: y_j \leq y}} P(X = x_i)P(Y = y_j). \end{aligned}$$

The double summation may contain infinite number of terms, but we can invoke Fubini theorem to write it as iterated summation (see also Example 23)

$$\begin{aligned} \sum_{\substack{i: x_i \leq x \\ j: y_j \leq y}} P(X = x_i)P(Y = y_j) &= \sum_{i: x_i \leq x} P(X = x_i) \sum_{j: y_j \leq y} P(Y = y_j) \\ &= P(X \leq x)P(Y \leq y). \end{aligned}$$

Hence

$$P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y), \quad \forall x, y \in \mathbb{R}.$$

By Theorem 61  $X$  and  $Y$  are independent.

2. Conversely suppose  $X$  and  $Y$  are independent. Then  $\forall i, j$ , the events  $\{X = x_i\}$ ,  $\{Y = y_j\}$  are independent, so (9.3) holds.  $\square$

In general, the probability of the event  $\{X \in (a, b], Y \in (c, d]\}$  can be expressed in terms of joint distribution function  $F$ ,

$$P(X \in (a, b], Y \in (c, d]) = F(b, d) - F(b, c) - F(a, d) + F(a, c).$$

**Thm 62.** *Suppose that the collection of events*

$$A_{IJ} = \{A_{ij} : i \in I, j \in J\}$$

*are independent. Here  $I, J$  are finite or infinite index sets. Let*

$$\mathcal{F}_i = \sigma(A_{ij} : j \in J), \forall i \in I.$$

*Then  $\mathcal{F}_1, \mathcal{F}_2, \dots$  are independent.*

**REMARK 4.** *Independence of infinite number of events is defined as any finite subcollection being independent.*

PROOF.  $\forall i \in I$ , denote

$$\mathcal{A}_i = \{\text{all finite intersections of } A_{i1}, A_{i2}, \dots\}.$$

Then  $\mathcal{A}_1, \mathcal{A}_2, \dots$  are  $\pi$ -systems and  $\mathcal{F}_i = \sigma(\mathcal{A}_i)$ . Hence the conclusion follows from Theorem 60.  $\square$

By inspecting the proof, we see that the above theorem extends to the case where the index set  $J$  varies with  $i \in I$ . As an application, we have the following useful result which states that functions of disjoint subgroups of independent random variables are independent.

For ease of writing, we introduce the notation of indexing by set, for example,

if  $I = \{i_1, \dots, i_l\}$ , then  $X_I$  means  $(X_{i_1}, \dots, X_{i_l})$ .

**Thm 63.** *Divide the independent random variables  $X_1, \dots, X_n$  into disjoint subgroups  $X_{I_1}, \dots, X_{I_k}$ , where  $I_1, I_2, \dots, I_k \subset \{1, 2, \dots, n\}$  are disjoint,  $\bigcup_{i=1}^k I_i = \{1, 2, \dots, n\}$ . If  $g_1(x_{I_1}), \dots, g_k(x_{I_k})$  are measurable functions, then  $g_1(X_{I_1}), \dots, g_k(X_{I_k})$  are independent.*

PROOF. Note  $\sigma(g_s(X_{I_s})) \subset \sigma(X_{I_s})$ ,  $s \in \{1, \dots, k\}$ , hence it is enough to show that  $\sigma(X_{I_1}), \dots, \sigma(X_{I_k})$  are independent. Each  $\sigma(X_{I_i})$  can be generated by  $\{\sigma(X_j) : j \in I_i\}$ . Now the proof is completed by applying Theorem 62 to

$$A_{IJ} = \{\sigma(X_j) : j \in I_i\}.$$

□

## 9.2. Independence and expectation.

**Thm 64.** *Suppose that  $X_1, \dots, X_n$  are independent with respective distribution  $\mu_i$ . Then  $(X_1, \dots, X_n)$  has the joint distribution  $\mu_1 \times \dots \times \mu_n$ .*

PROOF.  $\forall B_1, \dots, B_n \in \mathcal{B}(\mathbb{R})$ , using independence and the definition of product measure

$$\begin{aligned} P(X_1 \in B_1, \dots, X_n \in B_n) &= \prod_{i=1}^n P(X_i \in B_i) = \prod_{i=1}^n \mu_i(B_i) \\ &= \mu_1 \times \dots \times \mu_n(B_1 \times \dots \times B_n). \end{aligned}$$

The class of measurable rectangles

$$\{B_1 \times \cdots \times B_n : B_1, \dots, B_n \in \mathcal{B}(\mathbb{R})\}$$

is a  $\pi$ -system. Therefore by uniqueness (Theorem 10),  $\mu_1 \times \cdots \times \mu_n$  agrees with the joint distribution of  $(X_1, \dots, X_n)$ .  $\square$

**Thm 65.** *Suppose that  $X, Y$  are independent with respective distribution  $\mu$  and  $\nu$ ,  $h : \mathbb{R}^2 \mapsto \mathbb{R}$  is measurable. If  $h \geq 0$  or  $E|h(X, Y)| < \infty$ , then*

$$Eh(X, Y) = \int_{\mathbb{R}^2} h(x, y) d\mu(x) d\nu(y).$$

*In particular if  $h(x, y) = f(x)g(y)$ , then*

$$Ef(X)g(Y) = Ef(X) \cdot Eg(Y).$$

**PROOF.** By Theorem 64, the induced probability of  $(X, Y)$  is given by the product of  $\mu$  and  $\nu$ . Then using the general change of variables formula (Section 5.3, Formula 9.2), we can write

$$Eh(X, Y) = \int_{\mathbb{R}^2} h(x, y) d\mu(x) d\nu(y)$$



The remaining conclusion follows easily. □

Inductively using the above theorem we get the following expectation formula for independent random variables.

**Thm 66.** *Suppose that  $X_1, \dots, X_n$  are independent and either (a)  $X_i \geq 0, \forall i$  or (b)  $E|X_i| < \infty, \forall i$ . Then*

$$EX_1 \cdots X_n = EX_1 \cdots EX_n.$$

PROOF. The nonnegative case is immediate from Theorem 65. To prove case (b), applying the nonnegative case to  $|X_1|$  and  $|X_2|$ , we have

$$E|X_1 X_2| = E|X_1| \cdot E|X_2|.$$

Since the RHS is finite, so  $E|X_1 X_2| < \infty$  and by Theorem 65,

$$EX_1 X_2 = EX_1 \cdot EX_2.$$

Now the nonnegative case is again applicable to  $|X_1 X_2|$  and  $|X_3|$ ,

$$E|X_1 X_2 X_3| = E|X_1 X_2| \cdot E|X_3|.$$

So  $E|X_1X_2X_3| < \infty$  and Theorem 65 can be invoked to get

$$EX_1X_2X_3 = EX_1 \cdot EX_2 \cdot EX_3.$$

The procedure continues until  $X_n$  is processed, then the proof is completed.  $\square$

**Def 41.** *The covariance of  $X$  and  $Y$  is defined as*

$$\text{Cov}(X, Y) = E(X - EX)(Y - EY) = EXY - EX \cdot EY.$$

**Def 42.**  *$X$  and  $Y$  are uncorrelated if  $\text{Cov}(X, Y) = 0$ , i.e.,*

$$EXY = EX \cdot EY.$$

That  $X$  and  $Y$  are independent implies that they are uncorrelated, but not vice versa.

**Example 28.** *Suppose that  $X, Y$  are jointly distributed as below*

	$Y = -1$	$Y = 0$	$Y = 1$
$X = -1$	0	1/4	0
$X = 0$	1/4	0	1/4
$X = 1$	0	1/4	0

Then  $EX = EY = EXY = 0$ , but  $X, Y$  are not independent by Example 27, since

$$P(X = 0, Y = 0) = 0 \neq \frac{1}{4} = P(X = 0)P(Y = 0).$$

**Thm 67.** Suppose that  $X_1, \dots, X_n$  are pairwise uncorrelated and  $EX_i^2 < \infty$ . Then

$$\text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n).$$

PROOF. Denote  $X = X_1 + \dots + X_n$ . We have

$$\begin{aligned}\text{Var}(X) &= E(X - EX)^2 = E\left(\sum_{i=1}^n (X_i - EX_i)\right)^2 \\ &= E\left(\sum_{i=1}^n (X_i - EX_i)^2 + \sum_{i \neq j} (X_i - EX_i)(X_j - EX_j)\right) \\ &= \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n).\end{aligned}$$

We have used that  $X_1, \dots, X_n$  are pairwise uncorrelated, hence  $\forall i \neq j$ ,  

$$E((X_i - EX_i)(X_j - EX_j)) = EX_iX_j - EX_iEX_j = 0.$$

□

**9.3. Sum of independent random variables.** "Independent and identically distributed" is abbreviated as **i.i.d.**

**Example 29. (*Binomial distribution*)** Let  $p \in (0, 1)$ ,  $X_1, \dots, X_n$  i.i.d.  $\sim \text{Bernoulli}(p)$ . Define  $S_n = \sum_{i=1}^n X_i$ . Then

$$P(S_n = k) = C_n^k p^k (1 - p)^{n-k}, \quad k = 0, \dots, n.$$

The distribution of  $S_n$  is the binomial distribution, written  $S_n \sim \text{Bin}(n, p)$ . By employing linearity of integration, Example 14 and Theorem 67,

$$ES_n = np, \quad \text{Var}(S_n) = np(1 - p).$$

**Example 30 (The problem of points).** A coin with head probability  $p$  is flipped repeatedly. Gambler A wins one point if head appears

on a toss, otherwise gambler  $B$  wins one point. Whoever reaches first the finishing line wins the game. Suppose that gambler  $A$  and  $B$  are  $m$  and  $n$  points away from the finishing line. We intend to find the probability  $W(m, n)$  that gambler  $A$  wins the game. Imagine tossing the coins  $m + n - 1$  times, then gambler  $A$  wins the game if and only if heads show up at least  $m$  times, the probability is

$$W(m, n) = \sum_{k=m}^{m+n-1} C_{m+n-1}^k p^k (1-p)^{m+n-1-k}.$$

The probability of  $A$  winning the game can be categorized based on the outcome of the first toss. If the first toss is a head, the probability of  $A$  winning the game afterwards would be  $W(m-1, n)$ , otherwise  $W(m, n-1)$ . Therefore the recursion holds

$$W(m, n) = p \cdot W(m-1, n) + (1-p) \cdot W(m, n-1).$$

The equation may be solved by observing the boundary conditions,

$$W(0, j) = 1 \text{ for } j = n, n-1, \dots, 1$$

and

$$W(i, 0) = 0 \text{ for } i = m, m - 1, \dots, 1.$$

**Example 31 (Sum of Binomials).** Suppose that  $X \sim \text{Bin}(m, p)$  and  $Y \sim \text{Bin}(n, p)$  are independent, then  $X + Y \sim \text{Bin}(m + n, p)$ .

PROOF. For any  $k \in \{0, 1, \dots, m + n\}$ ,

$$\begin{aligned} P(X + Y = k) &= \sum_{i=0}^k P(X = i, Y = k - i) = \sum_{i=0}^k P(X = i)P(Y = k - i) \\ &= \sum_{i=0}^k C_m^i p^i (1 - p)^{m-i} \cdot C_n^{k-i} p^{k-i} (1 - p)^{n-k+i} \\ &= p^k (1 - p)^{m+n-k} \sum_{i=0}^k C_m^i C_n^{k-i} = C_{m+n}^k p^k (1 - p)^{m+n-k}. \end{aligned}$$

The last equality is due to Vandermonde identity. □

**Thm 68 (Convolution).** *Suppose that  $X, Y$  are independent with distribution functions  $F$  and  $G$ . Then*

$$P(X + Y \leq z) = \int_{\mathbb{R}} G(z - x) dF(x) = \int_{\mathbb{R}} F(z - y) dG(y)$$

Recall that  $dF, dG$  are notational variants for the corresponding Lebesgue-Stieltjes measure (Remark 3).

PROOF. Denote by  $\mu, \nu$  the distribution of  $X, Y$ . By Fubini theorem

$$\begin{aligned} P(X + Y \leq z) &= \int_{\mathbb{R}} \left( \int_{(-\infty, z-x]} d\nu(y) \right) d\mu(x) \\ &= \int_{\mathbb{R}} \left( \int_{(-\infty, z-y]} d\mu(x) \right) d\nu(y). \end{aligned}$$

The integrands respectively equal  $G(z - x)$  and  $F(z - y)$ . □

**Def 43.** *The random vector  $(X_1, \dots, X_n)$  has continuous distribution if there exists a function  $p \geq 0$  such that*

$$P((X_1, \dots, X_n) \in A) = \int_A p(x_1, \dots, x_n) dx_1 \cdots dx_n, \quad \forall A \in \mathcal{B}(\mathbb{R}^n).$$

*The function  $p$  is called the (joint) density of  $(X_1, \dots, X_n)$ .*

The definition can be equivalently<sup>3</sup> stated as: The random vector  $(X_1, \dots, X_n)$  has continuous distribution if there exists a function  $p \geq 0$  such that the joint distribution function  $F$  has

$$F(x_1, \dots, x_n) = \int_{(-\infty, x]} p(s_1, \dots, s_n) ds_1 \cdots ds_n, \quad \forall x = (x_1, \dots, x_n) \in \mathbb{R}^n,$$

where  $(-\infty, x] = (-\infty, x_1] \times \cdots \times (-\infty, x_n]$ .

In view of the definition, if  $X$  has density  $p$ , then its distribution function  $F$  has

$$F(x) = \int_{-\infty}^x p(s) ds, \quad \forall x.$$

---

<sup>3</sup>See previous section.



**Example 32 (Exponential distribution).**  $X$  has exponential distribution with parameter  $\lambda > 0$ , written  $X \sim \text{Exp}(\lambda)$ , if  $X$  has density

$$p(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Then  $EX = \lambda^{-1}$ ,  $\text{Var}(X) = \lambda^{-2}$ . In probability and statistics, the exponential distribution models the distribution of the waiting time before an event occurs.

**Example 33 (Normal distribution).**  $X$  has normal distribution with parameter  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ , written  $X \sim \mathcal{N}(\mu, \sigma^2)$ , if  $X$  has density

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}.$$

Then  $EX = \mu$ ,  $\text{Var}(X) = \sigma^2$ .

**Thm 69.** Suppose that  $X, Y$  are independent with distribution functions  $F$  and  $G$ . If  $X$  has density  $p_X$ , then  $Z = X + Y$  has density

$$h(z) = \int p_X(z - y)dG(y).$$

If also  $Y$  has density  $p_Y$ , then

$$h(z) = \int p_X(z - y)p_Y(y)dy.$$

PROOF. The convolution Theorem 68 now becomes

$$P(X + Y \leq z) = \int_{\mathbb{R}} \left( \int_{-\infty}^{z-y} p_X(x)dx \right) dG(y).$$

Combining a change of variable  $u = x + y$  with Fubini theorem, we obtain

$$\begin{aligned} P(X + Y \leq z) &= \int_{\mathbb{R}} \left( \int_{-\infty}^z p_X(u - y)du \right) dG(y) \\ &= \int_{-\infty}^z \left( \int_{\mathbb{R}} p_X(u - y)dG(y) \right) du \end{aligned}$$

□

**Example 34 (Sum of normal distributions).**  $X, Y$  are jointly normal, denoted by  $(X, Y) \sim \mathcal{N}(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$ , if the joint density  $p(x, y)$  is given by

$$\frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x-\mu_1}{\sigma_1}\right)\left(\frac{y-\mu_2}{\sigma_2}\right) + \left(\frac{y-\mu_2}{\sigma_2}\right)^2\right]\right\}.$$

Find the density of  $X, Y$  and  $Z = X + Y$ .

PROOF. Tedious calculations are omitted. The answers are

$$X \sim \mathcal{N}(\mu_1, \sigma_1^2), \quad Y \sim \mathcal{N}(\mu_2, \sigma_2^2).$$

and

$$Z \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + 2\rho\sigma_1\sigma_2 + \sigma_2^2).$$

□

**Def 44.** Gamma function is defined for  $\alpha > 0, \beta > 0$ ,

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx = \beta^\alpha \int_0^\infty x^{\alpha-1} e^{-\beta x} dx.$$

*Note the first integral does not depend on  $\beta$ .*

The following properties of Gamma functions are easy to verify that  $\Gamma(1) = 1$  and

$$\Gamma(\alpha + 1) = \alpha\Gamma(\alpha), \quad \forall \alpha > 0.$$

If  $n$  is a positive integer, then

$$\Gamma(n + 1) = n!.$$

**Def 45.** *Beta function is defined for  $\alpha > 0$ ,  $\beta > 0$ ,*

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1}(1-x)^{\beta-1}dx.$$

**LEMMA 70.** *Beta function is related to Gamma function by the equation*

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}.$$

PROOF. We start with from the definition of Gamma function,

$$\begin{aligned}\Gamma(\alpha)\Gamma(\beta) &= \int_0^\infty u^{\alpha-1}e^{-u}du \int_0^\infty v^{\beta-1}e^{-v}dv \\ &= \int_0^\infty \int_0^\infty u^{\alpha-1}v^{\beta-1}e^{-u-v}dudv.\end{aligned}$$

Now perform a change of variables,

$$u = st, \ v = s(1 - t), \text{ for } s > 0, \ 0 < t < 1.$$

The Jacobian determinant

$$\frac{\partial(u, v)}{\partial(s, t)} = \det \begin{pmatrix} t & s \\ 1 - t & -s \end{pmatrix} = -s.$$

Hence by Fubini theorem

$$\begin{aligned}\Gamma(\alpha)\Gamma(\beta) &= \int_0^1 \int_0^\infty s^{\alpha-1} t^{\alpha-1} s^{\beta-1} (1-t)^{\beta-1} e^{-s} ds dt \\ &= \int_0^\infty s^{\alpha+\beta-1} e^{-s} ds \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt \\ &= \Gamma(\alpha + \beta) B(\alpha, \beta).\end{aligned}$$

□

**Example 35 (Gamma distribution).** *A random variable follows a Gamma distribution with parameter  $\alpha > 0$ ,  $\beta > 0$ , if it has the density*

$$p(x) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, & \text{for } x > 0, \\ 0, & \text{for } x \leq 0. \end{cases}.$$

*Written Gamma( $\alpha, \beta$ ).*

Note Gamma(1,  $\beta$ ) is exponential distribution with parameter  $\beta$ , Exp( $\beta$ ).

**Thm 71 (Sum of Gamma).** *Suppose that  $X_i \sim \text{Gamma}(\alpha_i, \beta)$ ,  $i = 1, \dots, n$  are independent. Then*

$$Y = X_1 + \cdots + X_n \sim \text{Gamma}(\alpha_1 + \cdots + \alpha_n, \beta).$$

*If  $\alpha_i = 1$ , we can imagine  $n$  customers in a queue, each must wait time  $X_i$  for service once reaching the head of the queue. The average service rate is  $\beta$ . Then  $Y$  is the total waiting time of all  $n$  customers.*

**PROOF.** It suffices to prove for  $i = 2$ . Write  $p_{X_1}, p_{X_2}$  for the densities of  $X_1, X_2$ . Then the density of  $Y = X_1 + X_2$  is given by Theorem [69](#),

$$p_Y(y) = \int_0^y p_{X_1}(y-x)p_{X_2}(x)dx, \quad y > 0.$$

The integration is from 0 to  $y$ , since the densities  $p_{X_1}(x_1)$  and  $p_{X_2}(x_2)$  are non-zero only if  $x_1 > 0$ ,  $x_2 > 0$ . Hence

$$\begin{aligned}
p_Y(y) &= \int_0^y \frac{\beta^{\alpha_1}}{\Gamma(\alpha_1)}(y-x)^{\alpha_1-1}e^{-\beta(y-x)}\frac{\beta^{\alpha_2}}{\Gamma(\alpha_2)}x^{\alpha_2-1}e^{-\beta x}dx \\
&= \frac{\beta^{\alpha_1+\alpha_2}e^{-\beta y}}{\Gamma(\alpha_1)\Gamma(\alpha_2)}\int_0^y (y-x)^{\alpha_1-1}x^{\alpha_2-1}dx \\
&=_{(x=yt)} \frac{\beta^{\alpha_1+\alpha_2}e^{-\beta y}}{\Gamma(\alpha_1)\Gamma(\alpha_2)}\int_0^y (y-yt)^{\alpha_1-1}(yt)^{\alpha_2-1}ydt \\
&= \frac{\beta^{\alpha_1+\alpha_2}y^{\alpha_1+\alpha_2-1}e^{-\beta y}}{\Gamma(\alpha_1)\Gamma(\alpha_2)}\int_0^1 (1-t)^{\alpha_1-1}t^{\alpha_2-1}dt.
\end{aligned}$$

Using Lemma 70, we get

$$p_Y(y) = \frac{\beta^{\alpha_1+\alpha_2}}{\Gamma(\alpha_1+\alpha_2)}y^{\alpha_1+\alpha_2-1}e^{-\beta y} \sim \text{Gamma}(\alpha_1+\alpha_2, \beta).$$

□



**Def 46.** A random variable  $X$  has Beta distribution with parameter  $\alpha > 0, \beta > 0$ , written  $X \sim \text{Beta}(\alpha, \beta)$ , if it has the density

$$p(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 < x < 1,$$

where  $B(\alpha, \beta)$  is the Beta function. Lemma 70 confirms that the integration of  $p(x)$  over  $(0, 1)$  equals one.

**Example 36.** If  $X \sim \text{Beta}(\alpha, \beta)$  with  $\alpha > 0, \beta > 0$ , then using the properties of Gamma functions and Lemma 70, it is easy to check that

$$EX = \frac{\alpha}{\alpha + \beta}, \quad \text{Var}(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

**9.4. Construction of independent sequence.** Given a probability measure  $\mu$ , we intend to construct a sequence of independent random variables so that each has distribution  $\mu$ . Consider the probability space

$$(\Omega, \mathcal{F}, P) = ((0, 1], \mathcal{B}((0, 1]), \text{Lebesgue measure}).$$

Let  $X \sim U((0, 1])$  and take its **dyadic expansion**

$$(9.4) \quad X(\omega) = \sum_{k=1}^{\infty} \frac{X_k(\omega)}{2^k}.$$

The sequence  $X_k$  is determined by the algorithm:  $X_1 = 0$  if  $X \in (0, 1/2]$  (we can omit the left boundary point 0 since it carries zero Lebesgue measure), and 1 if  $X \in (1/2, 1]$ . Since  $X$  is distributed uniformly, it has equal probability of landing in  $(0, 1/2]$  or  $(1/2, 1]$ , hence the random variable  $X_1$  has Bernoulli distribution with parameter  $1/2$ ,

$$P(X_1 = 0) = P(X \in (0, 1/2]) = 1/2,$$

and

$$P(X_1 = 1) = P(X \in (1/2, 1]) = 1/2.$$

If  $X_1, \dots, X_{k-1}$  are already determined, then split into two halves the interval where  $X$  locates and define  $X_k = 0$  if  $X$  is on the left half, and 1 on the right half. As early,

$$X_k \sim \text{Bernoulli}(1/2).$$

Then we see by induction that for all  $\omega \in \Omega$ ,  $X(\omega)$  is bracketed by a sequence of intervals:  $X(\omega) \in D_n$  for all  $n \geq 1$  where  $D_n$  is the **dyadic interval of rank  $n$** ,

$$(9.5) \quad D_n = \left( \sum_{k=1}^n \frac{X_k(\omega)}{2^k}, \sum_{k=1}^n \frac{X_k(\omega)}{2^k} + \frac{1}{2^n} \right].$$

This implies that every dyadic expansion defined this way is **non-terminating**, otherwise if there is  $n_0$  such that  $X_n = 0$  for  $n > n_0$ , then

$$X(\omega) = \sum_{n=1}^{n_0} \frac{X_n(\omega)}{2^n}.$$

But this contradicts that  $X(\omega) \in D_{n_0}$ . A consequence of the property is the (natural) uniqueness of dyadic expansion. As an example, between the two mathematically equivalent expressions of  $1/2$ ,

$$\frac{1}{2} + \frac{0}{2} + \frac{0}{2} + \cdots \text{ and } \frac{0}{2} + \frac{1}{2} + \frac{1}{2} + \cdots,$$

the algorithm always chooses the one with infinitely many 1s, i.e., the former is *not* a dyadic expansion by the algorithm.

In numerics, the dyadic expansion of  $X \in (0, 1]$  is nothing but the binary representation

$$X = 0.X_1X_2X_3 \cdots$$

with  $X_1, X_2, \dots$  interpreted as binary digits.

LEMMA 72. *Every  $X \sim U((0, 1])$  has a unique (non-terminating) dyadic expansion, i.e., two expansions of  $X$  generated by the algorithm must necessarily have equal coefficients  $X_k$ s.*

PROOF. The uniqueness follows from the algorithm itself. Another way to prove it is to compare the coefficients of the expansion. Observe that if  $X$  has the dyadic expansion (9.4), then the coefficients  $X_1, X_2, \dots$  are uniquely determined as functions of  $X$ . For this, we define an operator similar to the floor function. To be consistent with our algorithm, the operator should map  $x \in (n, n + 1]$  to  $n$  for any

$n \in \{0, 1, \dots\}$ . Thus the function is

$$[[x]] = \text{ceil}(x) - 1 \text{ for } x > 0,$$

where  $\text{ceil}(x) = \inf\{n \in \mathbb{N} : n \geq x\}$  is the smallest integer no less than  $x$ . Now we can succinctly write  $X_1$  as a function of  $X$ . Multiply (9.4) by 2,

$$2X = 2\left(\frac{X_1}{2}\right) + 2\left(\frac{X_2}{2^2} + \frac{X_3}{2^3} + \dots\right).$$

Since  $X_k \in \{0, 1\}$  and the expansion generated by the algorithm is non-terminating, so the second term cannot be zero,

$$2\left(\frac{X_2}{2^2} + \frac{X_3}{2^3} + \dots\right) \in (0, 1].$$

Therefore

$$[[2X]] = 2\left(\frac{X_1}{2}\right) = X_1.$$

Similarly we have

$$[[2^2X]] = 2^2\left(\frac{X_1}{2} + \frac{X_2}{2^2}\right).$$

Hence

$$X_2 = 2^2\left(\frac{X_1}{2} + \frac{X_2}{2^2}\right) - 2 \cdot 2\left(\frac{X_1}{2}\right) = [[2^2X]] - 2[[2X]].$$

By induction that, for  $k \geq 1$ ,

$$X_k = [[2^kX]] - 2[[2^{k-1}X]].$$

Therefore  $X_1, X_2, \dots$  are uniquely determined by  $X$ . □

**REMARK 5.** *If the base space is  $\Omega = [0, 1)$  and  $X_1 = 0$  if  $X \in [0, 1/2)$ , and 1 if  $X \in [1/2, 1)$ , and so on, then the generated dyadic expansion of  $1/2$  would be*

$$\frac{1}{2} = \frac{1}{2} + \frac{0}{2} + \frac{0}{2} + \dots$$

*In particular the expansion coefficients  $X_1, X_2, \dots$  are uniquely determined as*

$$X_k = [2^k X] - 2[2^{k-1} X] \text{ for } k \geq 1.$$

*Recall  $[\cdot]$  is the floor function that takes out the integer part.*

LEMMA 73. *If  $X \sim U((0, 1])$  has the dyadic expansion (9.4), then  $X_1, X_2, \dots$  i.i.d  $\sim \text{Bernoulli}(1/2)$ .*

PROOF. That each  $X_k$  has Bernoulli distribution with parameter  $1/2$  is clearly from the algorithm. To show that  $X_1, X_2, \dots$  are independent, it is enough to verify that for any  $n \geq 1$ ,  $i_1, i_2, \dots, i_n \in \{0, 1\}$ ,

$$P(X_1 = i_1, \dots, X_n = i_n) = P(X_1 = i_1) \cdots P(X_n = i_n).$$

But this is immediate once we observe that the RHS equals  $1/2^n$  by construction and

$$\{X_1 = i_1, \dots, X_n = i_n\} = \left\{ X \in \left( \sum_{k=1}^n \frac{i_k}{2^k}, \sum_{k=1}^n \frac{i_k}{2^k} + \frac{1}{2^n} \right] \right\}.$$

Since  $X$  is uniform, the event on the RHS has probability  $1/2^n$ . Hence the desired independence follows.  $\square$

**Thm 74.** *If  $X \sim U((0, 1])$ , then there are i.i.d random variables  $X_1, X_2, \dots \sim \text{Bernoulli}(1/2)$  so that*

$$(9.6) \quad X = \sum_{k=1}^{\infty} \frac{X_k}{2^k}.$$

*Conversely, if (9.6) holds for independent  $X_1, X_2, \dots$  with distribution  $\text{Bernoulli}(1/2)$ , then  $X \sim U((0, 1])$ .*

**PROOF.** The first conclusion is contained in Lemma 73. Suppose that (9.6) holds for a sequence of independent Bernoulli random variables  $X_1, X_2, \dots$ , with parameter  $1/2$ . Define  $Y_n = \sum_{k=1}^n X_k/2^k$ . Clearly  $Y_n \rightarrow X$ , a.s., hence converges in distribution. Thus it suffices to show that the pointwise limit of the distribution functions of  $Y_n$  is indeed the distribution function of  $U((0, 1])$ , so that  $X \sim U((0, 1])$ . For each



$n \geq 1$ ,  $Y_n$  takes  $2^n$  distinct values,

$$Y_n \in \left\{ \frac{0}{2^n}, \frac{1}{2^n}, \dots, \frac{2^n - 1}{2^n} \right\},$$

so the distribution function of  $Y_n$  is a step function with jumps at  $i/2^n$ ,  $i = 0, 1, \dots, 2^n - 1$  and

$$(9.7) \quad P\left(Y_n = \frac{i}{2^n}\right) = \frac{1}{2^n}.$$

For any  $y \in (0, 1]$ , the dyadic expansion generation algorithm tells us that for all  $r \geq 1$ ,  $y$  is contained in a dyadic interval  $D_r$  of rank  $r$  (Imagine that  $y$  plays the role of  $\omega$  in (9.5)). The boundary points of  $D_r$  can be explicitly written down in terms of  $y$ . Each dyadic interval of rank  $r$  has length  $1/2^r$ , the number of these intervals that come before  $y$  is  $[y/(1/2^r)]$ , so

$$D_r = \left( \frac{1}{2^r} [2^r y], \frac{1}{2^r} [2^r y] + \frac{1}{2^r} \right].$$

It follows for all  $r \geq 1$ ,

$$P\left(Y_n \leq \frac{1}{2^r} [2^r y]\right) \leq P(Y_n \leq y) \leq P\left(Y_n \leq \frac{1}{2^r} [2^r y] + \frac{1}{2^r}\right).$$

Now we take  $r = n$  and show that both of the extreme terms converge to  $y$  as soon as  $n \rightarrow \infty$ , it would follow that  $P(Y_n \leq y)$  converges to the identity function  $y \mapsto y$  on  $(0, 1]$ , the distribution function of  $U((0, 1])$ . We only compute the LHS, the RHS is handled similarly. By (9.7),

$$p_n(y) \triangleq P\left(Y_n \leq \frac{1}{2^n} [2^n y]\right) = \sum_{i=0}^{[2^n y]} P\left(Y_n = \frac{i}{2^n}\right) = \frac{[2^n y] \wedge (2^n - 1)}{2^n},$$

Note if  $y = 1$ ,  $[2^n y] = 2^n$ , then the effective upper bound of the summation is  $2^n - 1$ , since the maximal value of  $Y_n$  is  $(2^n - 1)/2^n$ . Since

$$\frac{2^n y - 1}{2^n} < \frac{[2^n y]}{2^n} \leq \frac{2^n y}{2^n},$$

so  $p_n(y) \rightarrow y$ . The proof is completed. □

**Thm 75.** *Given a probability measure  $\mu$ , there exist a sequence of i.i.d random variables with  $\mu$  being the common distribution.*

PROOF. Let  $X(\omega) = \omega$ , then  $X \sim U((0, 1])$ . By Theorem 74, there are independent  $X_1, X_2, \dots \sim \text{Bernoulli}(1/2)$  so that (9.6) holds. Let  $\alpha$  be the one-to-one mapping from  $\mathbb{N} \times \mathbb{N}$  to  $\mathbb{N}$  and define

$$Y_{ij} = X_{\alpha(i,j)} \text{ and } U_i = \sum_{j=1}^{\infty} \frac{Y_{ij}}{2^j}, \quad i = 1, 2, \dots$$

Then  $U_1, U_2, \dots$  are independent by Theorem 62. By Theorem 74 again, we see that  $U_1, U_2, \dots$  have uniform distribution  $U((0, 1])$ . Let  $F$  be the distribution function associated with  $\mu$  and  $F^{-1}$  the inverse distribution function, then Theorem 63 and the proof of Theorem 23 tell us that  $F^{-1}(U_1), F^{-1}(U_2), \dots$  are i.i.d with common distribution  $F$ .  $\square$

## 10. 大数律 Law of large numbers

### 10.1. $L_2$ weak law.

**Thm 76** ( $L_2$  weak law). Suppose that  $X_1, \dots, X_n$  are pairwise uncorrelated with  $EX_i = \mu$  and  $\text{Var}(X_i) \leq C < \infty$ . Let  $S_n = X_1 + \dots + X_n$ . Then

$$\frac{S_n}{n} \rightarrow \mu \text{ in } L_2 \text{ and probability.}$$

PROOF. We only show  $L_2$  convergence, which will give convergence of probability via Markov inequality. Using the variance of sum formula (Theorem 67), we have

$$E \left| \frac{S_n}{n} - \mu \right|^2 = \text{Var} \left( \frac{S_n}{n} \right) = \frac{\text{Var}(S_n)}{n^2} = \frac{\sum_i \text{Var}(X_i)}{n^2} \leq \frac{C}{n} \rightarrow 0.$$

□

An important special case of the  $L_2$  weak law is the following.

**Thm 77.** Suppose that  $X_1, \dots, X_n$  are i.i.d with  $EX_i = \mu$  and  $\text{Var}(X_i) = \sigma^2$ . Let  $S_n = X_1 + \dots + X_n$ . Then

$$\frac{S_n}{n} \rightarrow \mu \text{ in } L_2 \text{ and probability.}$$

Below is a probabilistic proof of Weierstrass approximation theorem.

**Example 37 (Bernstein polynomial).** *f is continuous on  $[0, 1]$ . Define*

$$f_n(x) = \sum_{k=0}^n f\left(\frac{k}{n}\right) C_n^k x^k (1-x)^{n-k}, \quad \forall x \in [0, 1].$$

*Then*

$$\sup_{x \in [0, 1]} |f_n(x) - f(x)| \rightarrow 0.$$

PROOF. Observe that, if we let  $S_n \sim \text{Bin}(n, x)$ , then

$$f_n(x) = Ef(S_n/n).$$

Hence

$$\begin{aligned} (10.1) \quad |f_n(x) - f(x)| &= |Ef(S_n/n) - f(x)| = |E[f(S_n/n) - f(x)]| \\ &\leq E|f(S_n/n) - f(x)|. \end{aligned}$$

$\forall \varepsilon > 0$ , we want to bound the rightmost expectation in terms of  $\varepsilon$ . Since  $f$  is continuous on  $[0, 1]$  and hence uniformly continuous, we can fix  $\delta$  small so that

$$|f(s) - f(t)| < \varepsilon \text{ if } |s - t| < \delta.$$

Let  $M = \sup_{x \in [0, 1]} f(x)$  and  $A_n = \{\omega : |S_n(\omega)/n - x| < \delta\}$ . Then (10.1) continues

$$\begin{aligned} |f_n(x) - f(x)| &\leq E(|f(S_n/n) - f(x)|1_{A_n}) + E(|f(S_n/n) - f(x)|1_{A_n^c}) \\ &\leq \varepsilon + 2MP(|S_n/n - x| \geq \delta) \\ &\leq_{(e_1)} \varepsilon + 2M \frac{\text{Var}(S_n/n)}{\delta^2} =_{(e_2)} \varepsilon + 2M \frac{x(1-x)}{n\delta^2} \\ &\leq \varepsilon + \frac{M}{2n\delta^2} \leq 2\varepsilon, \end{aligned}$$

as soon as  $n$  is large so that  $M/(2n\delta^2) \leq \varepsilon$ , where  $(e_1)$  uses Markov inequality, and  $(e_2)$  Example 29.  $\square$

LEMMA 78. If  $b_n$  satisfies  $\text{Var}(S_n)/b_n^2 \rightarrow 0$ , then

$$\frac{S_n - ES_n}{b_n} \rightarrow 0 \text{ in } L_2 \text{ and probability.}$$

PROOF. We have

$$\text{Var}\left(\frac{S_n - ES_n}{b_n}\right) = \frac{\text{Var}(S_n)}{b_n^2} \rightarrow 0.$$

□

**Example 38 (Coupon collector's problem).** Suppose there are  $n$  types of coupons. You get one coupon each time you open a box of candy, and the coupon is equally likely to be any of the  $n$  types. We are interested in the time  $T_n$  to collect a complete set of coupons. Let  $\tau_0^n = 0$  and

$\tau_k^n =$  the first time we have  $k$  different coupons,  $k = 1, \dots, n$ .

Then

$$T_n = \tau_n^n = \sum_{k=1}^n (\tau_k^n - \tau_{k-1}^n).$$

It is readily seen that the waiting times  $\{\tau_k^n - \tau_{k-1}^n\}_{k=1}^n$  between two types of coupons are independent and each has geometric distribution,

$$\tau_k^n - \tau_{k-1}^n \sim \text{Geom}\left(1 - \frac{k-1}{n}\right).$$

Example 16 tells us that

$$ET_n = \sum_{k=1}^n \left(1 - \frac{k-1}{n}\right)^{-1} = n \sum_{m=1}^n m^{-1} \approx n \log n.$$

and

$$\text{Var}(T_n) \leq \sum_{k=1}^n \left(1 - \frac{k-1}{n}\right)^{-2} = n^2 \sum_{m=1}^n m^{-2} \leq n^2 \sum_{m=1}^{\infty} m^{-2}.$$

Since  $\sum_{m=1}^{\infty} m^{-2}$  is convergent, if we take  $b_n = n \log n$ , then

$$\frac{\text{Var}(T_n)}{b_n^2} \leq \frac{n^2 \sum_{m=1}^{\infty} m^{-2}}{(n \log n)^2} \rightarrow 0.$$



So Lemma 78 gives

$$\frac{T_n - n \sum_{m=1}^n m^{-1}}{n \log n} \rightarrow 0 \text{ in probability.}$$

It follows that

$$\frac{T_n}{n \log n} \rightarrow 1 \text{ in probability.}$$

This tells us that  $T_n$  is roughly  $n \log n$ .

**Example 39 (Random permutation).** A permutation of  $\{1, \dots, n\}$  is a one-to-one mapping from  $\{1, \dots, n\}$  to itself. There are  $n!$  permutations in total. We are interested in the expected number of cycles in a randomly chosen permutation. As an example, we look at the permutation  $i \mapsto \pi(i)$ ,

$$\begin{array}{rcccccc} i : & 1 & 2 & 3 & 4 & 5 & 6 \\ \pi(i) : & 2 & 5 & 6 & 4 & 1 & 3 \end{array}$$

Starting with 1, we follow the route of mapping

$$1 \rightarrow \pi(1) \rightarrow \pi^2(1) \rightarrow \pi^3(1) \rightarrow \dots$$

Since  $\pi^3(1) = 1$ , we get a cycle

$$1 \rightarrow 2 \rightarrow 5 \rightarrow 1.$$

We use brackets to indicate cycles, so we have the first cycle (125), and the remaining cycles are (36), (4). The original permutation can now be simply written as the decomposition

$$(125)(36)(4).$$

The representations of a permutation as a mapping and block decomposition are equivalent. This is indeed how random permutation generation algorithm works: decompose  $1, \dots, n$  into disjoint blocks, at the  $k$ -th position of the decomposition, the algorithm has choices with equal probability among the  $n - k$  numbers that have not been seen so far (if the block is to grow) plus the first number of the current block the algorithm is in (if the block is to close so that a cycle is formed).

Therefore, if we define the indicator random variables

$$X_{n,k} = \begin{cases} 1, & \text{a closing bracket occurs after the } k\text{-th} \\ & \text{position in the decomposition,} \\ 0, & \text{otherwise.} \end{cases}$$

then

$$S_n \triangleq X_{n,1} + \cdots + X_{n,n}$$

gives the total number of cycles in the permutation and

$$P(X_{n,k} = 1) = \frac{1}{n - k + 1}.$$

It can also be verified that for  $k < l$ ,

$$P(X_{n,k} = 1, X_{n,l} = 1) = \frac{1}{n - k + 1} \cdot \frac{1}{n - l + 1},$$

which implies that  $X_{n,k}$ ,  $X_{n,l}$  are independent (by Example 27 and Theorem 58). The same direct computation can prove the independence

of  $\{X_{n,1}, \dots, X_{n,n}\}$ . Then

$$ES_n = \sum_{k=1}^n EX_{n,k} = \sum_{k=1}^n k^{-1}$$

and noting  $X_{n,k}^2 = X_{n,k}$  we have

$$\text{Var}(S_n) \leq \sum_{k=1}^n EX_{n,k}^2 = \sum_{k=1}^n EX_{n,k} = \sum_{k=1}^n k^{-1}.$$

Now applying Lemma 78 with  $b_n = (\log n)^{0.5+\varepsilon}$ ,  $\varepsilon > 0$ ,

$$\frac{S_n - \sum_{k=1}^n k^{-1}}{b_n} \rightarrow 0 \text{ in probability.}$$

It follows that, if  $\varepsilon = 0.5$ ,

$$\frac{S_n}{(\log n)^{0.5+\varepsilon}} \rightarrow 1 \text{ in probability.}$$

The arbitrariness of  $\varepsilon$  indicates that  $(\log n)^{0.5}$  is a threshold for the convergence.

Random permutation is commonly used in applications from coding to games, one example is the **100 prisoners riddle**. 100 prisoners, who are numbered from 1 to 100, are offered a last chance to be pardoned. At a room, there is a cupboard with 100 drawers. 100 numbers from 1 to 100 are randomly put into these drawers. The prisoners enter the room one by one. Each prisoner can open up to 50 drawers. No communications are allowed. If every prisoner finds their numbers, all prisoners are set free, otherwise all will be sentenced. If every prisoner randomly opens 50 drawers, the survival probability would be  $(1/2)^{100}$ . The prisoners need to figure out the best strategy to follow.

The numbers in the drawers form a permutation  $\pi$  of  $\{1, \dots, 100\}$ , the drawer labelled with  $i$  contains the number  $\pi(i)$ . The permutation is decomposed as collections of cycles. The strategy is thus to enter the correct cycle containing the wanted number. For the prisoner with number  $i_0$ , the first drawer to open is the one labelled with  $i_0$ , subsequently with label  $\pi(i_0)$ ,  $\pi^2(i_0)$ , ... Since every number is in some cycle, there is  $k$ ,  $1 \leq k \leq n$ , so that  $\pi^k(i_0) = i_0$ , i.e., the wanted number  $i_0$  would be found after opening the drawer labelled with  $\pi^{k-1}(i_0)$ . The

prisoners survive the test if the random permutation in the drawer contains no cycle of length strictly greater than 50 (there is at most one in every permutation). The probability of a random permutation containing a cycle of length  $k$  is

$$\frac{C_{100}^k \cdot (k-1)! \cdot (100-k)!}{100!}.$$

Therefore the survival probability of all prisoners is then equal to

$$1 - \sum_{k=51}^{100} \frac{C_{100}^k \cdot (k-1)! \cdot (100-k)!}{100!} = 1 - \sum_{k=51}^{100} \frac{1}{k}.$$

## 10.2. Weak law of large numbers.

**Thm 79** (Weak law for triangular arrays). *Consider the triangular array of random variables  $X_{n,k}$ ,  $k = 1, \dots, n$ ,*

$$\begin{array}{ccccccc} & & X_{1,1} & & & & \\ & & & & & & \\ & X_{2,1} & & X_{2,2} & & & \\ & \dots & & & & & \\ X_{n,1} & & \dots & & X_{n,k} & \dots & X_{n,n} \end{array}$$

*Random variables in each row are pairwise independent. Let  $b_n > 0$  satisfies*

(10.2)

$$(i) \sum_{k=1}^n P(|X_{n,k}| > b_n) \rightarrow 0; \quad (ii) \frac{\sum_{k=1}^n \text{Var}\left(X_{n,k} 1_{|X_{n,k}| \leq b_n}\right)}{b_n^2} \rightarrow 0.$$

*If we set  $S_n = \sum_{k=1}^n X_{n,k}$ ,  $a_n = \sum_{k=1}^n E\left(X_{n,k} 1_{|X_{n,k}| \leq b_n}\right)$ , then*

$$\frac{S_n - a_n}{b_n} \rightarrow 0 \text{ in probability.}$$

PROOF. Let

$$\bar{S}_n = \sum_{k=1}^n X_{n,k} 1_{|X_{n,k}| \leq b_n} \text{ and } Z_n = \frac{S_n - a_n}{b_n}.$$

We have  $\forall \varepsilon > 0$ ,

$$\begin{aligned} P(|Z_n| > \varepsilon) &= P(|Z_n| > \varepsilon, S_n \neq \bar{S}_n) + P(|Z_n| > \varepsilon, S_n = \bar{S}_n) \\ &\leq P(S_n \neq \bar{S}_n) + P(|Z_n| > \varepsilon, S_n = \bar{S}_n). \end{aligned}$$

By assumption (i),

$$P(S_n \neq \bar{S}_n) \leq \sum_{k=1}^n P(|X_{n,k}| > b_n) \rightarrow 0.$$



Now using assumption (ii) and that  $X_{n,i}, X_{n,j}, i \neq j$  are independent, we have

$$\begin{aligned} P(|Z_n| > \varepsilon, S_n = \bar{S}_n) &\leq P\left(\left|\frac{\bar{S}_n - a_n}{b_n}\right| > \varepsilon\right) \leq \frac{\text{Var}(\bar{S}_n)}{\varepsilon^2 b_n^2} \\ &= \frac{\sum_{k=1}^n \text{Var}(X_{n,k} 1_{|X_{n,k}| \leq b_n})}{\varepsilon^2 b_n^2} \rightarrow 0. \end{aligned}$$

The proof is complete. □

LEMMA 80. *If  $X \geq 0$ ,  $\varphi$  differentiable with  $\varphi' > 0$  and  $\varphi(0) = 0$ , then*

$$\int \varphi(X) dP = \int_0^\infty \varphi'(t) P(X > t) dt.$$

PROOF. An application of Theorem 56 with  $Y = \varphi(X)$  gives,

$$\begin{aligned}\int Y dP &= \int_0^\infty P(Y > s) ds \\ &=_{(e_1)} \int_0^\infty \varphi'(t) P(Y > \varphi(t)) dt = \int_0^\infty \varphi'(t) P(X > t) dt.\end{aligned}$$

We have performed in  $(e_1)$  a change of variable  $s = s(t) = \int_0^t \varphi'$ .  $\square$

**Thm 81 (Weak law of large numbers).** *Let  $X_1, \dots, X_n$  be i.i.d with*

$$(10.3) \quad xP(|X_1| > x) \rightarrow 0 \text{ as } x \rightarrow 0.$$

*If we set  $S_n = X_1 + \dots + X_n$ ,  $\mu_n = E(X_1 1_{|X_1| \leq n})$ , then*

$$\frac{S_n}{n} - \mu_n \rightarrow 0 \text{ in probability.}$$

PROOF. We want to apply Theorem 79 with  $X_{n,k} = X_k$  and  $b_n = n$ . To do this, we need to verify condition (10.2). First note that

$$\sum_{k=1}^n P(|X_{n,k}| > b_n) = nP(|X_1| > n) \rightarrow 0$$

and

$$\frac{\sum_{k=1}^n \text{Var}\left(X_{n,k} 1_{|X_{n,k}| \leq b_n}\right)}{b_n^2} = \frac{\text{Var}\left(X_1 1_{|X_1| \leq n}\right)}{n} \leq \frac{E\left(\left(X_1 1_{|X_1| \leq n}\right)^2\right)}{n},$$

recalling Example 17. Thus the proof would be completed if we show that

$$\frac{E\left(\left(X_1 1_{|X_1| \leq n}\right)^2\right)}{n} \rightarrow 0.$$

By Lemma 80,

$$(10.4) \quad E\left(\left(X_1 1_{|X_1| \leq n}\right)^2\right) = \int_0^\infty 2tP(|X_1| 1_{|X_1| \leq n} > t) dt.$$

Note the integrand has

$$P(|X_1|1_{|X_1| \leq n} > t) = P(|X_1| > t, |X_1| \leq n),$$

which gives the expression,

$$P(|X_1|1_{|X_1| \leq n} > t) = \begin{cases} P(|X_1| > t) - P(|X_1| > n), & t < n \\ 0, & t \geq n \end{cases}$$

Hence upon substituting the above in (10.4) we obtain

$$\frac{E\left((X_1 1_{|X_1| \leq n})^2\right)}{n} \leq \frac{1}{n} \int_0^n 2tP(|X_1| > t)dt.$$

Using the assumption that  $tP(|X_1| > t) \rightarrow 0$  as  $t \rightarrow \infty$ , we see that the RHS converges to zero, which completes the proof.  $\square$

A sufficient condition for (10.3) is  $E|X_1| < \infty$ . Indeed, by dominated convergence theorem,

$$xP(|X_1| > x) \leq E(|X_1|1_{|X_1| > x}) \rightarrow 0 \text{ as } x \rightarrow \infty.$$

So  $E|X_1| < \infty$  implies (10.3), and is thus a stronger condition, but the latter is not much weaker since by Lemma 80, for  $0 < \varepsilon < 1$ ,

$$\begin{aligned} E|X|^{1-\varepsilon} &= \int_0^\infty (1-\varepsilon)t^{-\varepsilon}P(X > t)dt \\ &= \int_0^1 (1-\varepsilon)t^{-\varepsilon}P(X > t)dt + \int_1^\infty (1-\varepsilon)t^{-\varepsilon}P(X > t)dt \\ &\leq \int_0^1 t^{-\varepsilon}dt + \int_1^\infty t^{-(1+\varepsilon)}tP(X > t)dt < \infty. \end{aligned}$$

**Thm 82.** *Let  $X_1, \dots, X_n$  be i.i.d with  $E|X_1| < \infty$ . If we set  $S_n = X_1 + \dots + X_n$ , then*

$$\frac{S_n}{n} \rightarrow EX_1 \text{ in probability.}$$

**PROOF.** Let  $\mu = EX_1$ ,  $\mu_n = E(X_1 1_{|X_1| \leq n})$ . As we have already seen that  $E|X_1| < \infty$  implies (10.3), so we can employ Theorem 81 to

conclude that,  $\forall \varepsilon > 0$ ,

$$P\left(\left|\frac{S_n}{n} - \mu_n\right| > \varepsilon\right) \rightarrow 0.$$

Since  $\mu_n \rightarrow \mu$  by dominated convergence theorem, we have  $|\mu_n - \mu| < \varepsilon$  for large  $n$ , therefore

$$P\left(\left|\frac{S_n}{n} - \mu\right| > 2\varepsilon\right) \leq P\left(\left|\frac{S_n}{n} - \mu_n\right| > \varepsilon\right).$$

It follows that  $S_n/n - \mu \rightarrow 0$  in probability.  $\square$

In the example below, we will see that weak law can exist even if the condition of Theorem 82 fails:  $E|X_1| = \infty$ .

**Example 40 (St. Petersburg paradox).** *A single-player game begins with an initial wager of 2 dollars and a fair coin. The coin is tossed repeatedly. Each time a tail comes up, the wager is doubled. The game ends if head appears. So if the first toss is head, the game ends and the player receives 2 dollars. The expected amount the player*

would receive is

$$2 \cdot \frac{1}{2} + 2^2 \cdot \frac{1}{2^2} + 2^3 \cdot \frac{1}{2^3} + \cdots = \infty.$$

Paradoxically, no one would pay an infinite amount to play a game. We want to use the weak law Theorem 79 to find the right value of the game. The idea is to see where the average game value goes after playing several rounds of the game. Let  $X_1, X_2, \dots$  be independent with values in  $\{2^m : m = 1, 2, \dots\}$  and satisfy

$$P(X_k = 2^m) = 2^{-m}.$$

To apply Theorem 79, we need to find  $b_n > 0$  so that (10.2) is satisfied with  $X_{n,k} = X_k$ . Since

$$nP(X_1 > b_n) = n \sum_{m: 2^m > b_n} 2^{-m}$$

and

$$nb_n^{-2} \text{Var}(X_1 1_{X_1 \leq b_n}) \leq nb_n^{-2} E(X_1 1_{X_1 \leq b_n})^2 \leq nb_n^{-2} \sum_{m: 2^m \leq b_n} 2^{2m} \cdot 2^{-m}$$

So (10.2)(i) and (ii) translate as requiring

$$(10.5) \quad n \sum_{m:2^m > b_n} 2^{-m} \rightarrow 0 \text{ and } nb_n^{-2} \sum_{m:2^m \leq b_n} 2^{2m} \cdot 2^{-m} \rightarrow 0.$$

We assume that

$$m(n) \triangleq \log_2 b_n = \log_2 n + K(n),$$

where  $K(n)$  is chosen so that  $m(n)$  is an integer. Then

$$n \sum_{m:2^m > b_n} 2^{-m} \leq n2^{-m(n)} = 2^{-K(n)}$$

and

$$\begin{aligned} nb_n^{-2} \sum_{m:2^m \leq b_n} 2^{2m} \cdot 2^{-m} &= nb_n^{-2} \frac{2(2^{m(n)+1} - 1)}{2 - 1} \\ &\leq 4nb_n^{-2} 2^{m(n)} \leq 4nb_n^{-1} = 4 \cdot 2^{-K(n)}. \end{aligned}$$



Hence all it takes for (10.5) to hold is  $K(n) \rightarrow \infty$  while keeping  $m(n)$  an integer. Thus Theorem 79 tells us that, with  $S_n = \sum_{k=1}^n X_k$ ,

$$(10.6) \quad \frac{S_n - a_n}{n2^{K(n)}} \rightarrow 0 \text{ in probability.}$$

where

$$a_n = nE(X_1 1_{|X_1| \leq b_n}) = n \sum_{m: 2^m \leq b_n} 2^m \cdot 2^{-m} = nm(n) = n(\log_2 n + K(n)).$$

To draw a meaningful conclusion, we choose  $K(n)$  so that

$$\frac{a_n}{n2^{K(n)}} = \frac{\log_2 n + K(n)}{2^{K(n)}} \rightarrow 1.$$

In particular, if  $K(n) \approx \log_2 \log_2 n$  for large  $n$ , then the above is satisfied and (10.6) gives

$$\frac{S_n}{n \log_2 n} \rightarrow 1 \text{ in probability.}$$

*This says that the average  $S_n/n$  of  $n$  rounds of the game is close to  $\log_2 n$ , which should therefore be a reasonable price for the game.*

### 10.3. Borel-Cantelli lemma and applications.

LEMMA 83 (**The first Borel-Cantelli lemma**). *We have*

$$\sum_{k=1}^{\infty} P(A_k) < \infty \text{ implies } P(A_n \text{ i.o.}) = 0.$$

PROOF. We have

$$P\left(\limsup_n A_n\right) = \lim_n P\left(\bigcup_{k=n}^{\infty} A_k\right) \leq \lim_n \sum_{k=n}^{\infty} P(A_k) = 0.$$

□

The next is a typical application of Borel-Cantelli lemma, the application to the strong law of large numbers is postponed to Theorem 88.

LEMMA 84. Suppose that  $\varepsilon_n \geq 0$  satisfies  $\sum_n \varepsilon_n < \infty$  and the random variables  $X_n$  have

$$\sum_{n=1}^{\infty} P(|X_{n+1} - X_n| > \varepsilon_n) < \infty.$$

Then there exists a finite random variable  $X$  so that  $X_n \rightarrow X$ , a.s.

PROOF. Let  $A_n = \{|X_{n+1} - X_n| > \varepsilon_n\}$  and  $A^* = \limsup_n A_n$ . Then

for  $\omega \in (A^*)^c = \bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} A_k^c$  if and only if there is  $m(\omega)$  satisfying

$$|X_{n+1}(\omega) - X_n(\omega)| \leq \varepsilon_n \text{ for } n \geq m(\omega),$$

hence  $\{X_n(\omega)\}$  is Cauchy and converges to some finite limit, say  $X^*(\omega)$ . Define  $X = 0$  for  $\omega \in A^*$  and  $X = X^*$  otherwise. Then  $X$  is a random variable since  $A^*$  is measurable. Using the first Borel-Cantelli lemma, we have  $P(A^*) = 0$  which shows that  $X_n \rightarrow X$ , a.s.  $\square$

LEMMA 85 (**The second Borel-Cantelli lemma**). *If  $A_1, \dots, A_n$  are independent, then*

$$\sum_{n=1}^{\infty} P(A_n) = \infty \text{ implies } P(A_n \text{ i.o.}) = 1.$$

PROOF. It suffices to show  $P\left(\bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} A_k^c\right) = 0$ , which clearly follows if we can show that  $P\left(\bigcap_{k=n}^{\infty} A_k^c\right) = 0$  for all  $n$ . By independence and  $1 - x \leq e^{-x}$ ,

$$P\left(\bigcap_{k=n}^N A_k^c\right) = \prod_{k=n}^N (1 - P(A_k)) \leq \prod_{k=n}^N e^{-P(A_k)} = \exp\left\{-\sum_{k=n}^N P(A_k)\right\}.$$

The latter converges to zero as  $N \rightarrow \infty$ . Hence

$$P\left(\bigcap_{k=n}^{\infty} A_k^c\right) = \lim_N P\left(\bigcap_{k=n}^N A_k^c\right) = 0.$$

□

**COROLLARY 4 (Zero-One law).** *If  $A_1, \dots, A_n$  are independent, then  $P(A_n \text{ i.o.}) = 0$  or 1 according as  $\sum_{n=1}^{\infty} P(A_n)$  converges or diverges.*

Borel-Cantelli lemmas are easier to understand when translated into the language of random variables. Let

$$(10.7) \quad S_n = \sum_{k=1}^n 1_{A_k} \text{ and } S = \sum_{k=1}^{\infty} 1_{A_k},$$

then we have the translation of Borel-Cantelli lemmas in the language of random variables,

Lemma 83 :  $ES < \infty$  implies  $S < \infty$  *a.s.*

Lemma 85 : If  $A_1, A_2, \dots$  are independent, then  
 $ES = \infty$  implies  $S = \infty$  *a.s.*

Note by monotone convergence theorem  $\lim_n ES_n = ES$ .

With the random variable translation, we can easily show that the second Borel-Cantelli lemma (Lemma 85) continues to hold if independence is replaced with pairwise independence.

LEMMA 86 (**The second Borel-Cantelli lemma**). *If  $A_1, \dots, A_n$  are pairwise independent, then*

$$\sum_{n=1}^{\infty} P(A_n) = \infty \text{ implies } P(A_n \text{ i.o.}) = 1.$$

PROOF. Let  $S_n = \sum_{k=1}^n 1_{A_k}$  and  $S = \sum_{k=1}^{\infty} 1_{A_k}$ . We see that the desired conclusion is equivalent to  $P(S < \infty) = 0$ . By pairwise independence,

$$\text{Var}(S_n) = \sum_{k=1}^n \text{Var}(1_{A_k}) \leq \sum_{k=1}^n P(A_k) = ES_n.$$

Since  $S_n \leq S$ , we have

$$\begin{aligned} P(S < ES_n/2) &\leq P(S_n < ES_n/2) \\ &\leq P(|S_n - ES_n| > ES_n/2) \leq \frac{4\text{Var}(S_n)}{(ES_n)^2} \leq \frac{4}{ES_n}. \end{aligned}$$

Noting that  $ES_n \rightarrow \infty$  by assumption, we proceed to write

$$P(S < \infty) = \lim_n P(S < ES_n/2) \leq \lim_n \frac{4}{ES_n} = 0.$$

□

REMARK 6. Assume the same conditions as Lemma 86. A slight modification of the proof of Lemma 86 yields that,  $\forall \delta > 0$ ,

$$(10.8) \quad P(|S_n - ES_n| > \delta ES_n) \leq \frac{\text{Var}(S_n)}{\delta^2 (ES_n)^2} \leq \frac{1}{\delta^2 ES_n}$$

So we have

$$ES_n \rightarrow \infty \text{ implies } \frac{S_n}{ES_n} \rightarrow 1 \text{ in probability.}$$

This can also be derived directly from Lemma 78.

Through a useful technique which we call the **method of subsequence**, we show that the above convergence can be strengthened and it is indeed almost sure. We will again see the use of the method of subsequence in the proof of the strong law (Theorem 89).

**Thm 87.** If  $A_1, \dots, A_n$  are pairwise independent, then

$$\sum_{n=1}^{\infty} P(A_n) = \infty \text{ implies } \frac{\sum_{k=1}^n 1_{A_k}}{\sum_{k=1}^n P(A_k)} \rightarrow 1, \text{ a.s.}$$



PROOF. Let  $N_n = \sum_{k=1}^n 1_{A_k}$ , we want to prove that  $N_n/EN_n \rightarrow 1$ , *a.s.* To do this, we proceed in two steps.

1. First we show that the conclusion is true for a subsequence. Let

$$\tau_k = \inf\{n : EN_n \geq k\} \text{ and } S_k = N_{\tau_k}.$$

By the definition  $EN_{\tau_{k^2}-1} < k^2$  (otherwise  $\tau_{k^2}$  would be no greater than  $\tau_{k^2} - 1$ , a contradiction). So

$$(10.9) \quad k^2 \leq ES_{k^2} = EN_{\tau_{k^2}-1} + E1_{A_{\tau_{k^2}}} < k^2 + 1.$$

Notice that the assumption  $EN_n \rightarrow \infty$  ensures that  $\tau_k < \infty$  and  $\tau_k \rightarrow \infty$  so that (10.9) holds for all  $k \geq 1$ , otherwise if  $\sup_n EN_n < \infty$ , then for large  $k$ ,  $\tau_k = \infty$  ( $\inf \emptyset = \infty$ ) and  $k^2 \leq ES_{k^2}$  would not hold. Now applying (10.8) to the subsequence  $k \mapsto S_{k^2}$  gives,  $\forall \delta > 0$ ,

$$P(|S_{k^2} - ES_{k^2}| > \delta ES_{k^2}) \leq \frac{1}{\delta^2 ES_{k^2}} \leq \frac{1}{\delta^2 k^2}.$$

So the first Borel-Cantelli Lemma (Lemma 83) shows that

$$\frac{S_{k^2}}{ES_{k^2}} \rightarrow 1, a.s.$$

**2.** Next we extend the conclusion from  $S_{k^2}$  to the whole sequence. Note that for any  $n$  with  $\tau_{k^2} \leq n \leq \tau_{(k+1)^2}$ ,

$$\frac{S_{k^2}}{ES_{(k+1)^2}} \leq \frac{N_n}{EN_n} \leq \frac{S_{(k+1)^2}}{ES_{k^2}},$$

which can be rewritten as

$$\frac{S_{k^2}}{ES_{k^2}} \cdot \frac{ES_{k^2}}{ES_{(k+1)^2}} \leq \frac{N_n}{EN_n} \leq \frac{S_{(k+1)^2}}{ES_{(k+1)^2}} \cdot \frac{ES_{(k+1)^2}}{ES_{k^2}}.$$

Thus the desired conclusion follows if  $ES_{(k+1)^2}/ES_{k^2} \rightarrow 1$ , but this is immediate since by the definition of  $S_{k^2}$ ,  $S_{(k+1)^2}$  and (10.9),

$$1 \leq \frac{ES_{(k+1)^2}}{ES_{k^2}} < \frac{(k+1)^2 + 1}{k^2} \rightarrow 1.$$

Therefore the proof is completed. □

**Example 41 (Record values).** Suppose that i.i.d random variables  $X_1, X_2, \dots$  from a continuous distribution function  $F$  are observed sequentially. Denote by  $A_k = \{X_k > X_i \text{ for } i = 1, \dots, k-1\}$  the event that a record occurs at the  $k$ -th random variable. We want to determine the asymptotics of the count

$$R_n = \sum_{k=1}^n 1_{A_k}$$

of record events in the first  $n$  random variables. Since the distribution function is continuous, the values of  $X_1, X_2, \dots, X_n$  are almost surely distinct<sup>4</sup>. By rearranging  $X_1, X_2, \dots, X_n$  in decreasing order, we obtain a permutation  $\pi_n$  over  $1, \dots, n$ , where all  $n!$  permutations are equally likely. The event  $A_k$  occurs if and only if the  $k$ -th position is the greatest among the first  $k$ , this is, in the language of permutation,  $\pi_k(k) = 1$ . Note that the permutation after the  $k$ -th position does not affect that of the first  $k$ . There are only one way to put the greatest (of

---

<sup>4</sup>Durrett 5th Exercise 2.1.5

the first  $k$ ) at the  $k$ -th position and the remaining can be permuted in any of  $(k-1)!$  ways. Hence

$$P(A_k) = P(\pi_k(k) = 1) = \frac{1 \cdot (k-1)!}{k!} = \frac{1}{k}.$$

The same idea generalizes to multiple record events, for example, for  $k < l$ ,

$$P(A_k A_l) = \frac{1 \cdot (l-1)!}{l!} \cdot \frac{1 \cdot (k-1)!}{k!} = \frac{1}{k} \cdot \frac{1}{l} = P(A_k)P(A_l).$$

With these it can be verified that  $A_1, A_2, \dots, A_n$  are independent. Now we can employ Theorem 87 to conclude that

$$\frac{R_n}{\log n} \rightarrow 1, \text{ a.s.}$$

Note the conclusion is independent of  $F$  as long as it is continuous.

**10.4. Strong law of large numbers.** Our first version of strong law of large numbers is a typical application of the first Borel-Cantelli lemma (Lemma 83)

**Thm 88.** Let  $X_1, \dots, X_n$  be i.i.d with  $EX_1^4 < \infty$ . If we set  $S_n = X_1 + \dots + X_n$ , then

$$\frac{S_n}{n} \rightarrow EX_1 \text{ a.s.}$$

PROOF. Assuming without loss of generality that  $EX_1 = 0$ , we observe that the desired conclusion amounts to,  $\forall \varepsilon > 0$ ,

$$(10.10) \quad P(|S_n| > n\varepsilon \text{ i.o.}) = 0.$$

We have by Markov inequality that

$$(10.11) \quad P(|S_n| > n\varepsilon) \leq \frac{ES_n^4}{(n\varepsilon)^4}.$$

Now

$$ES_n^4 = E\left(\sum_{1 \leq i, j, k, l \leq n} X_i X_j X_k X_l\right) = \sum_{1 \leq i, j, k, l \leq n} E(X_i X_j X_k X_l).$$

By the i.i.d assumption and that  $EX_1 = 0$ , we see from Theorem 66 that  $E(X_i X_j X_k X_l)$  is zero unless it is of either one of the form  $EX_i^4$ ,

$EX_i^2 X_j^2$  with  $i \neq j$ . There are respectively  $n$  and  $C_4^2 \cdot C_n^2 = 3n(n-1)$  of these terms (for the latter, pick two indices out of  $i, j, k, l$  and then two distinct random variables out of  $X_1, \dots, X_n$ ). Hence

$$ES_n^4 = nEX_i^4 + 3n(n-1)EX_i^2 X_j^2 = nEX_1^4 + 3n(n-1)(EX_1^2)^2 \leq Cn^2,$$

where  $C$  is a constant independent of  $n$ . Plugging this into (10.11), we obtain

$$P(|S_n| > n\varepsilon) \leq \frac{C}{n^2\varepsilon^4}.$$

Hence  $\sum_{n=1}^{\infty} P(|S_n| > n\varepsilon) < \infty$ , so (10.10) follows from the first Borel-Cantelli lemma (Lemma 83). □

The i.i.d and fourth order moment assumption of Theorem 88 can be weakened. Next we give Etemadi's proof of **Kolmogorov's strong law of large numbers** under pairwise independence and finite first order moment.

**Thm 89 (Strong law of large numbers).** *Suppose that  $X_1, \dots, X_n$  are pairwise independent identically distributed with  $E|X_1| < \infty$ . If we set  $S_n = X_1 + \dots + X_n$ , then*

$$\frac{S_n}{n} \rightarrow EX_1 \text{ a.s.}$$

PROOF. We start by observing that if the theorem holds for non-negative random variable, then

$$\frac{S_n}{n} = \frac{1}{n} \left( \sum_{k=1}^n X_k^+ - \sum_{k=1}^n X_k^- \right) \rightarrow EX_1^+ - EX_1^- = EX_1 \text{ a.s.}$$

So we can assume from now on that  $X_k \geq 0$ ,  $k \geq 1$ . As in Theorem 79, we define the truncated partial sum

$$\bar{S}_n = \sum_{k=1}^n X_k 1_{X_k \leq k}.$$

Let  $\alpha > 1$  and  $\tau_n = [\alpha^n]$ .

1. We first show that

$$\sum_{n=1}^{\infty} P(|\bar{S}_{\tau_n} - E\bar{S}_{\tau_n}| > \varepsilon \tau_n) < \infty.$$

As usual

$$\text{Var}(\bar{S}_{\tau_n}) = \sum_{k=1}^{\tau_n} \text{Var}(X_k 1_{X_k \leq k}) \leq \sum_{k=1}^{\tau_n} E(X_k^2 1_{X_k \leq k}) \leq \tau_n E(X_1^2 1_{X_1 \leq \tau_n}).$$

Hence

$$\sum_{n=1}^{\infty} P(|\bar{S}_{\tau_n} - E\bar{S}_{\tau_n}| > \varepsilon \tau_n) \leq \sum_{n=1}^{\infty} \frac{\text{Var}(\bar{S}_{\tau_n})}{\varepsilon^2 \tau_n^2} \leq \frac{1}{\varepsilon^2} E \left[ X_1^2 \sum_{n=1}^{\infty} \frac{1_{X_1 \leq \tau_n}}{\tau_n} \right].$$

For  $x = X_1(\omega) > 0$ , let  $n_x = \min\{n \in \mathbb{N} : \tau_n \geq x\}$ . By the definition we have  $\tau_n \geq \alpha^n/2$ <sup>5</sup> and  $\alpha^{n_x} \geq \tau_{n_x} \geq x$ , it follows that

$$\sum_{n=1}^{\infty} \frac{1_{X_1 \leq \tau_n}}{\tau_n} = \sum_{n \geq n_x} \frac{1}{\tau_n} \leq 2 \sum_{n \geq n_x} \alpha^{-n} = \frac{2\alpha^{-n_x}}{1 - \alpha^{-1}} \leq \frac{2x^{-1}}{1 - \alpha^{-1}} = \frac{2X_1^{-1}}{1 - \alpha^{-1}}.$$

---

<sup>5</sup>For  $z \geq 1$ ,  $z/2 < [z]$ : if  $z \in [1, 2)$ ,  $z/2 < 1 = [z]$ ; if  $z \geq 2$ ,  $z - [z] < z/2$ .



Therefore

$$\sum_{n=1}^{\infty} P(|\bar{S}_{\tau_n} - E\bar{S}_{\tau_n}| > \varepsilon \tau_n) \leq \frac{2}{\varepsilon^2(1 - \alpha^{-1})} EX_1 < \infty.$$

2. Next we claim that

$$\frac{S_{\tau_n}}{\tau_n} \rightarrow EX_1, \text{ a.s.}$$

With what we already have from step 1, we can invoke the first Borel-Cantelli lemma to obtain that

$$\frac{\bar{S}_{\tau_n} - E\bar{S}_{\tau_n}}{\tau_n} \rightarrow 0, \text{ a.s.}$$

But  $EX_k 1_{X_k \leq k} \rightarrow EX_1$  by dominated convergence theorem, it follows that  $E\bar{S}_{\tau_n}/\tau_n \rightarrow EX_1$ . Hence  $\bar{S}_{\tau_n}/\tau_n \rightarrow EX_1, \text{ a.s.}$  Since

$$\begin{aligned} \sum_{k=1}^{\infty} P(X_k 1_{X_k \leq k} \neq X_k) &\leq \sum_{k=1}^{\infty} P(X_k > k) \leq \int_0^{\infty} P(X_1 > t) dt \\ &= EX_1 < \infty, \end{aligned}$$

invoking again the first Borel-Cantelli lemma we get that

$$P(X_k 1_{X_k \leq k} \neq X_k \text{ i.o.}) = 0,$$

hence  $(S_{\tau_n} - \bar{S}_{\tau_n})/\tau_n \rightarrow 0$ , *a.s.* It follows that

$$\frac{S_{\tau_n}}{\tau_n} = \frac{S_{\tau_n} - \bar{S}_{\tau_n}}{\tau_n} + \frac{\bar{S}_{\tau_n}}{\tau_n} \rightarrow EX_1, \text{ a.s.}$$

**3.** Finally we conclude via the use of subsequence method. For any  $k$  satisfying  $\tau_n \leq k \leq \tau_{n+1}$ , since  $X_k \geq 0$ , we have

$$\frac{S_{\tau_n}}{\tau_{n+1}} \leq \frac{S_k}{k} \leq \frac{S_{\tau_{n+1}}}{\tau_n},$$

which we rewrite as

$$\frac{S_{\tau_n}}{\tau_n} \cdot \frac{\tau_n}{\tau_{n+1}} \leq \frac{S_k}{k} \leq \frac{S_{\tau_{n+1}}}{\tau_{n+1}} \cdot \frac{\tau_{n+1}}{\tau_n}.$$

But by the definition  $\tau_{n+1}/\tau_n \rightarrow \alpha$ , so it follows from step **2** that

$$\frac{1}{\alpha} EX_1 \leq \liminf \frac{S_k}{k} \leq \limsup \frac{S_k}{k} \leq \alpha EX_1, \text{ a.s.}$$

The proof is completed by sending  $\alpha \rightarrow 1$ . □

The next theorem shows that for i.i.d sequence, finite first moment  $E|X_1| < \infty$  in Theorem 89 is not only sufficient but also necessary for the strong law to hold.

**Thm 90.** *Let  $X_1, \dots, X_n$  be i.i.d with  $E|X_1| = \infty$ , then*

$$P(|X_n| > n \text{ i.o.}) = 1$$

*and setting  $S_n = X_1 + \dots + X_n$ ,*

$$P\left(\lim_n \frac{S_n}{n} \text{ exists and is finite}\right) = 0.$$

PROOF. We have by Lemma 80

$$E|X_1| = \int_0^\infty P(|X_1| > t) dt \leq \sum_{n=0}^\infty P(|X_1| > n)$$

Since  $E|X_1| = \infty$ , we infer from the second Borel-Cantelli lemma that  $P(|X_n| > n \text{ i.o.}) = 1$ . To prove the remaining conclusion, let

$$C = \left\{ \omega : \lim_n \frac{S_n(\omega)}{n} \text{ exists and is finite} \right\}.$$

We claim that  $C$  does not intersect  $\{|X_n| > n \text{ i.o.}\}$ , thus has  $P(C) = 0$ . If  $\omega \in C \cap \{|X_n| > n \text{ i.o.}\}$ , then

$$\begin{aligned} \left| \frac{S_{n+1}(\omega)}{n+1} - \frac{S_n(\omega)}{n} \right| &= \left| \frac{S_n(\omega)}{n+1} - \frac{S_n(\omega)}{n} + \frac{X_{n+1}(\omega)}{n+1} \right| \\ &\geq \left| \frac{X_n(\omega)}{n+1} \right| - \left| \frac{S_n(\omega)}{n+1} - \frac{S_n(\omega)}{n} \right|. \end{aligned}$$

Whence

$$\limsup_n \left| \frac{S_{n+1}(\omega)}{n+1} - \frac{S_n(\omega)}{n} \right| \geq 1.$$

But this contradicts  $\omega \in C$ , therefore  $C \cap \{|X_n| > n \text{ i.o.}\}$  must be empty.  $\square$

The strong law of large numbers holds whenever  $EX_1$  exists in the extended sense, i.e., at least one of  $EX_1^+$ ,  $EX_1^-$  is finite.

**Thm 91.** *Let  $X_1, \dots, X_n$  be i.i.d with  $EX_1^+ = \infty$ ,  $EX_1^- < \infty$ . Set  $S_n = X_1 + \dots + X_n$ . Then*

$$\frac{S_n}{n} \rightarrow \infty \text{ a.s.}$$

PROOF. Define  $X_k^M = X_k \wedge M$ ,  $S_n^M = X_1^M + \dots + X_n^M$ , then by the assumption,  $X_1^M, \dots, X_n^M$  are i.i.d with finite expectation. Hence Theorem 89 tells us that for  $M > 0$ ,

$$\liminf_n \frac{S_n}{n} \geq \lim_n \frac{S_n^M}{n} \rightarrow EX_1^M, \text{ a.s.}$$

Now let  $M \rightarrow \infty$ . □

The following is an application of the strong law to statistics, particularly we will prove the Glivenko-Cantelli theorem which is usually referred to as the fundamental theorem of statistics.

**Example 42 (Empirical distribution function).** *Let  $X_1, \dots, X_n$  be i.i.d with distribution function  $F$ . Fix  $x \in \mathbb{R}$ . We can estimate the value  $F(x)$  as below. Define*

$$F_n(x, \omega) = \frac{\sum_{k=1}^n 1_{\{X_k \leq x\}}(\omega)}{n}.$$

*The mapping  $x \mapsto F_n(x, \omega)$  is the so called empirical distribution function. By the strong law (Theorem 89), for each  $x$ , there is an exception set  $A_x$  of zero probability,*

$$\lim_n F_n(x, \omega) = E1_{\{X_1 \leq x\}} = F(x), \quad \omega \in A_x^c.$$

*But the theorem below says more: the exception set  $A_x$  can be independent of  $x$ .*

In the following, we will drop  $\omega$  and simply write  $F_n$  for empirical distribution function. But keep in mind  $F_n$  is a random function.

Before we state the main theorem, it is worthwhile recalling the exercise <sup>6</sup>: if  $G_n, G$  are nondecreasing functions and  $G$  is bounded and

---

<sup>6</sup>Pku Textbook Chapter 2 Exercise 32

continuous, then  $G_n$  converges to  $G$  uniformly. But if  $G$  has discontinuities, then it is not immediately obvious that the convergence is uniform.

**Thm 92 (The Glivenko-Cantelli theorem).** *Let  $X_1, \dots, X_n$  be i.i.d with distribution function  $F$ . Then the associated empirical distribution function  $F_n$  has*

$$\sup_{x \in \mathbb{R}} |F_n(x, \omega) - F(x)|, \quad \omega\text{-a.s.}$$

**PROOF. 1.** As with Example 42, an application of the strong law shows that, for each  $x$ , there is an exception set  $B_x$  of zero probability so that

$$F_n(x-, \omega) \triangleq \frac{\sum_{k=1}^n 1_{\{X_k < x\}}(\omega)}{n} \rightarrow F(x-), \quad \omega \in B_x^c.$$

Let  $m \geq 1$  and  $F^{-1}(z)$ ,  $z \in (0, 1)$ , be the inverse distribution function (see (4.1)). Define

$$x_{m,k} = \begin{cases} F^{-1}(k/m), & 1 \leq k \leq m-1, \\ -\infty, & k = 0, \\ +\infty, & k = m. \end{cases}$$

Since

$$F(F^{-1}(z)-) \leq z \leq F(F^{-1}(z)), \quad \forall z \in (0, 1),$$

we have for  $2 \leq k \leq m-1$ ,

$$F(x_{m,k}-) - F(x_{m,k-1}) \leq k/m - (k-1)/m = 1/m.$$

The inequality remains true for  $k = 1$  or  $m$  with the understanding that  $F(x_{m,0}) = 0$  and  $F(x_{m,m}-) = 1$ . For  $1 \leq k \leq m$  and  $x_{m,k-1} \leq x < x_{m,k}$ , we get by monotonicity,

$$\begin{aligned} F_n(x) - F(x) &\leq F_n(x_{m,k}-) - F(x_{m,k-1}) \\ &\leq F(x_{m,k}-) - F(x_{m,k-1}) + |F_n(x_{m,k}-) - F(x_{m,k}-)| \\ &\leq 1/m + |F_n(x_{m,k}-) - F(x_{m,k}-)|. \end{aligned}$$



Similarly,

$$\begin{aligned}
F(x) - F_n(x) &\leq F(x_{m,k}-) - F_n(x_{m,k-1}) \\
&\leq F(x_{m,k}-) - F(x_{m,k-1}) + |F(x_{m,k-1}) - F_n(x_{m,k-1})| \\
&\leq 1/m + |F(x_{m,k-1}) - F_n(x_{m,k-1})|
\end{aligned}$$

Note  $|F(x_{m,0}) - F_n(x_{m,0})| = 0$ ,  $|F_n(x_{m,m}-) - F(x_{m,m}-)| = 0$ .

**2.** Let  $\varepsilon > 0$ . For  $1 \leq k \leq m$ , denote by  $A_{m,k}$  the exception set for the convergence of  $F_n(x_{m,k-1})$  ( $A_{m,0} \triangleq \emptyset$ ), and  $B_{m,k}$  for the convergence of  $F_n(x_{m,k}-)$  ( $B_{m,m} \triangleq \emptyset$ ). Also let  $E = \bigcup_{m=1}^{\infty} \bigcup_{k=1}^m (A_{m,k} \cup B_{m,k})$ , then  $P(E) = 0$ . If

$$D_n(m) \triangleq \frac{1}{m} + \max_{1 \leq k \leq m} \{|F_n(x_{m,k}-) - F(x_{m,k}-)|, |F(x_{m,k-1}) - F_n(x_{m,k-1})|\},$$

then for any  $m$  satisfying  $1/m < \varepsilon/2$ , there is  $N_m$  such that

$$D_n(m) \leq \varepsilon, \quad n \geq N_m, \quad \omega \in E^c.$$

Any  $x \in \mathbb{R}$  is contained in some interval  $[x_{m,k-1}, x_{m,k})$ , so by step 1, as soon as  $n \geq N_m$ ,

$$|F(x) - F_n(x)| \leq D_n(m) \leq \varepsilon, \quad \omega \in E^c,$$

completing the proof. □

Another application of the strong law, now to renewal theory.

**Thm 93 (Renewal theory).** *Imagine a number of lightbulbs produced by the same manufacturer are available. At time 0, a lightbulb is lit up and replaced by a new one when it burns out. The lifetime of these lightbulbs are modeled by i.i.d random variables  $X_1, X_2, \dots$  with  $0 < EX_1 = \lambda^{-1} \leq \infty$ . We are interested in the number of lightbulbs that have burned out by time  $t > 0$ ,*

$$N_t = \sup\{n : T_n \leq t\} = \sum_{n=1}^{\infty} 1_{T_n \leq t},$$

where  $T_n = X_1 + \cdots + X_n$ . Then as  $t \rightarrow \infty$ ,

$$\frac{N_t}{t} \rightarrow \lambda, \text{ a.s.}$$

If additionally,  $X_1$  (hence every  $X_k$ ) does not concentrate on a single point, i.e. for all  $a \geq 0$ ,  $P(X_1 = a) \neq 1$ , then we obtain the elementary limit theorem,

$$\frac{EN_t}{t} \rightarrow \lambda.$$

PROOF. 1. For  $t > 0$ , by the definition,  $T_{N_t} \leq t < T_{N_t+1}$ , so

$$(10.12) \quad \frac{N_t + 1}{T_{N_t+1}} \cdot \frac{N_t}{N_t + 1} = \frac{N_t}{T_{N_t+1}} < \frac{N_t}{t} \leq \frac{N_t}{T_{N_t}}.$$

Since  $X_1, X_2, \dots$  are i.i.d with finite expectation  $EX_1 = \lambda^{-1}$ , we get from the strong law Theorem 89 that

$$(10.13) \quad \frac{T_n}{n} \rightarrow \lambda^{-1}, \text{ a.s.}$$

Moreover, we infer from the finite expectation assumption that outside an exception set,  $T_n < \infty$  for all  $n$ . This together with (10.13) implies

that almost surely  $T_n \rightarrow \infty$  as  $n \rightarrow \infty$ , hence almost surely

$$N_t < \infty \text{ for all } t > 0.$$

We claim that almost surely  $N_t \rightarrow \infty$  as  $t \rightarrow \infty$ . If this is not true for  $\omega$  in a set of positive probability, then  $\{N_t(\omega) : t > 0\}$  is bounded by some finite  $n(\omega)$ , this implies that  $T_n(\omega) > t$  for all  $t$  whenever  $n > n(\omega)$ . Sending  $t \rightarrow \infty$  would contradict that almost surely  $T_n < \infty$  for all  $n$ . Now almost surely, both of the extreme terms of (10.12) are well-defined and converge to  $\lambda$ . Thus  $N_t/t \rightarrow \lambda$  almost surely.

**2.** To prove the second conclusion, it suffices to show that  $N_t/t$ ,  $t \geq 1$ , is uniformly integrable (by Theorem 45). Note we are only interested in  $t \rightarrow \infty$ , so we have excluded  $0 < t < 1$ . Let

$$X_{k,a} = a1_{X_k > a} \text{ for } a > 0, \quad k = 1, 2, \dots$$

This may be thought of as ignoring the lightbulb whose lifetime is no greater than a threshold  $a$ , while longer lifetime is simply counted as  $a$ . As a result, the renewal time of these "modified lightbulbs" can only happen at the times that are multiples of  $a$ ,  $\{na : n \in \mathbb{N}\}$ . So each renewal can be regarded as a geometric random variable waiting for

the event  $\{X_k > a\}$  to occur. To make the geometric random variable meaningful, we require that  $a > 0$  satisfy

$$0 < p \triangleq P(X_1 > a) < 1.$$

We can do this since  $X_1$  does not concentrate on any single point. The frequency of the modified renewal is higher than its unmodified counterpart. So, if

$$N_{t,a} = \sup\{n : T_{n,a} \leq t\},$$

where  $T_{n,a} = X_{1,a} + \cdots + X_{n,a}$ , then  $N_t \leq N_{t,a}$ . Since  $N_{t,a}$  is less than the sum of  $[t/a]$  independent geometric random variables with parameter  $p$ , hence

$$EN_{t,a}^2 \leq E\left(\sum_{i=1}^{[t/a]} \text{Geom}(p)\right)^2 \leq c(1 + t + t^2),$$

where  $c$  is a constant depending on  $a$  and  $p$  only. Now for  $y > 0$ ,

$$P\left(\frac{N_t}{t} > y\right) \leq P\left(\frac{N_{t,a}}{t} > y\right) \leq \frac{EN_{t,a}^2}{y^2 t^2} \leq \frac{3c}{y^2}, \quad \forall t \geq 1.$$

Denote  $Y_t = N_t/t$ . Using Lemma 80, we obtain that

$$\begin{aligned} EY_t 1_{Y_t > y} &= \int_0^\infty P(Y_t 1_{Y_t > y} > s) ds \\ &= \int_0^y P(Y_t > y) ds + \int_y^\infty P(Y_t > s) ds \\ &\leq \frac{3c}{y} + \int_y^\infty \frac{3c}{s^2} ds \rightarrow 0 \text{ as soon as } y \rightarrow \infty. \end{aligned}$$

Whence  $N_t/t$ ,  $t \geq 1$ , is uniformly integrable. □

**REMARK 7.** *Another approach to the second conclusion of Theorem 93 is to modify the lightbulbs by cutoff,*

$$X_k^M = X_k \wedge M \text{ for } M > 0, \quad k = 1, 2, \dots$$

*The proof will rely on Wald's equation, but not need the assumption that the lifetime is not a fixed value. However the assumption that  $EX_1 > 0$  is still in place so that  $P(X_1 = 0) = 1$  does not happen.*

REMARK 8. *A particular application of the renewal theory assumes that the lifetime of a lightbulb is exponentially distributed with parameter  $\lambda$ . So on average, during a time period of length  $t$ , the number  $N_t$  of lightbulbs that burns out by time  $t$  approximates  $\lambda t$ . In this sense, the parameter of exponential distribution is interpreted as the rate of events, i.e., number of events per unit time.*

## 11. 中心极限定理 Central limit theorem

**11.1. Introduction.** The central limit theorem and the law of large numbers answer two seemingly disparate but related questions. Let  $X_1, \dots, X_n$  be i.i.d with  $E|X_1| < \infty$ ,  $S_n = X_1 + \dots + X_n$ . Then the strong law says that

$$\left| \frac{S_n}{n} - EX_1 \right| \rightarrow 0, \text{ a.s.}$$

Now the question is how large should  $n$  be so that  $|S_n/n - EX_1|$  is less than a given tolerance. This concerns the speed of convergence. An answer has already been hinted by the proof of the  $L_2$  weak law (Theorem 76) which goes as this

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

This means  $E(S_n/n - EX_1)^2$  grows at the same speed as  $n^{-1}$ , or roughly  $|S_n/n - EX_1|$  grows at the speed  $n^{-1/2}$ . The statement of course is very vague, but the central limit theorem will tell us more.



**11.2. From Poisson distribution to Stirling formula.** Before we start, we will review a basic limit theorem and on the way provide an intuitive proof of Stirling formula.

Suppose  $\lambda > 0$  is the number of events in a unit time interval, thus  $\lambda$  is the rate of events. The unit time interval is divided into  $n$  subintervals. Assume that *the probability of an event is approximated by the rate of the event and the probability of two events occurring in a small interval is negligible*. Hence if  $n$  is large, the probability of an event occurring on a subinterval can well be thought of as  $\lambda/n$ , and the total number of events that occur approximately follows the binomial distribution  $Bin(n, \lambda/n)$ . The following limit theorem makes this precise. Since the occurrence of an event on a subinterval of small size is rare, the theorem is commonly referred to as the **law of small numbers**.

**Thm 94 (Poisson approximation to Binomial).** *Let  $p = p(n) \rightarrow 0$ ,  $n \rightarrow \infty$  so that  $np \rightarrow \lambda > 0$ . Then for  $0 \leq k \leq n$ ,*

$$\lim_{n \rightarrow \infty} C_n^k p^k (1-p)^{n-k} = \frac{\lambda^k}{k!} e^{-\lambda}.$$

PROOF. Write

$$C_n^k p^k (1-p)^{n-k} = \frac{n(n-1) \cdots (n-k+1)}{k!} p^k (1-p)^{n-k}.$$

Since

$$p = \frac{\lambda}{n} + o\left(\frac{1}{n}\right) = \frac{1}{n}(\lambda + o(1)),$$

we have

$$n(n-1) \cdots (n-k+1) p^k = \frac{n(n-1) \cdots (n-k+1)}{n^k} (\lambda + o(1))^k \rightarrow 1,$$

and note  $k/n \rightarrow 0$ , so

$$(1-p)^{n-k} = \left[1 - \left(\frac{\lambda}{n} + o\left(\frac{1}{n}\right)\right)\right]^{n-k} \rightarrow e^{-\lambda}.$$

Therefore the conclusion is proved. □

The limit theorem indicates that Poisson distribution has the same bell shape as binomial distribution.

**Thm 95.** Let  $X \sim \text{Poisson}(\lambda)$  with  $\lambda > 0$ . Write  $p_k = P(X = k)$ ,  $k \geq 0$ . Then

(1) if  $\lambda$  is an integer,

$$p_0 < \cdots < p_{\lambda-1} = p_\lambda > p_{\lambda+1} \cdots$$

(2) if  $\lambda$  is not an integer,

$$p_0 < \cdots < p_{[\lambda]} > p_{[\lambda]+1} \cdots$$

PROOF. Compute  $p_{k+1}/p_k = \lambda/(k+1)$ . □

**Thm 96 (Stirling formula).**

$$n! \sim n^n e^{-n} \sqrt{2\pi n}$$

AN INTUITIVE PROOF. Let  $X \sim \text{Poisson}(\lambda)$ . Recall that  $EX = \lambda$ ,  $\text{Var}(X) = \lambda$ . The probability mass function of  $X$  is bell-shaped and peaks at  $\lambda$  (if it is an integer) which is close to the normal density with mean  $\lambda$  and variance  $\lambda$ , at least near the peak. So

$$\frac{\lambda^\lambda}{\lambda!} e^{-\lambda} \approx \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{(\lambda-\lambda)^2}{2\lambda}} = \frac{1}{\sqrt{2\pi\lambda}},$$

which gives

$$\lambda! \approx \lambda^\lambda e^{-\lambda} \sqrt{2\pi\lambda}.$$

□

**11.3. De Moivre-Laplace limit theorem.** Let  $X_1, \dots, X_n$  be i.i.d with

$$P(X_1 = -1) = P(X_1 = 1) = \frac{1}{2}.$$

Then  $EX_1 = 0$ ,  $\text{Var}(X_1) = 1$ . Let  $S_n = X_1 + \dots + X_n$ . We intend to show that, through proper standardization, the distribution of  $S_{2n}$  is close to the standard normal distribution  $\mathcal{N}(0, 1)$ . Since

$$ES_{2n} = 0, \quad \text{Var}(S_{2n}) = 2n,$$

in the spirit of standardization we should compute the distribution of  $S_{2n}/\sqrt{2n}$ . Note  $S_{2n}$  can only take even numbers between  $-2n$  and  $2n$ . So we need to compute the probability of the form

$$P\left(\frac{S_{2n}}{\sqrt{2n}} = \frac{2k}{\sqrt{2n}}\right) = P(S_{2n} = 2k).$$

For each  $2k$  in  $S_{2n}$ 's range (assuming  $2k \geq 0$ , the nonpositive case is symmetric), denote by  $u$  the number of  $X_k$ s that are  $+1$ , and  $d$  the number of  $X_k$ s that are  $-1$ , then

$$u - d = 2k, \quad u + d = 2n.$$

Hence

$$u = n + k, \quad d = n - k.$$

So

$$P(S_{2n} = 2k) = C_{2n}^{n+k} \cdot \frac{1}{2^{n+k}} \cdot \frac{1}{2^{n-k}} = \frac{(2n)!}{(n+k)!(n-k)!} \cdot \frac{1}{2^{2n}}.$$

Now using Stirling formula,  $P(S_{2n} = 2k)$  approximates

$$\begin{aligned} & \frac{(2n)^{2n} e^{-2n} \sqrt{2\pi(2n)}}{(n+k)^{n+k} e^{-n-k} \sqrt{2\pi(n+k)} \cdot (n-k)^{n-k} e^{-n+k} \sqrt{2\pi(n-k)}} \cdot \frac{1}{2^{2n}} \\ &= \frac{n^{2n}}{(n+k)^{n+k} (n-k)^{n-k}} \cdot \frac{\sqrt{2\pi(2n)}}{\sqrt{2\pi(n+k)} \sqrt{2\pi(n-k)}} \end{aligned}$$

which we rewrite as

$$\left(1 - \frac{k^2}{n^2}\right)^{-n} \left(1 + \frac{k}{n}\right)^{-k} \left(1 - \frac{k}{n}\right)^k \cdot \frac{1}{\sqrt{\pi n}} \frac{\sqrt{n^2}}{\sqrt{(n+k)(n-k)}}$$

It follows that, if  $2k/\sqrt{2n} \approx x$ , then  $k^2/n \approx -x^2/2$ , the above product asymptotically equals

$$e^{x^2/2} e^{-x^2/2} e^{-x^2/2} \cdot \frac{1}{\sqrt{\pi n}}.$$

Therefore we have

LEMMA 97. *Let  $X_1, \dots, X_n$  be i.i.d with*

$$P(X_1 = -1) = P(X_1 = 1) = \frac{1}{2}.$$

*Let  $S_n = X_1 + \dots + X_n$ . If  $2k/\sqrt{2n} \rightarrow x$ , then*

$$P(S_{2n} = 2k) \sim \frac{1}{\sqrt{\pi n}} e^{-\frac{x^2}{2}}.$$

Now we are in a position to compute, for  $a < b$ ,

$$\begin{aligned} P\left(a \leq \frac{S_{2n}}{\sqrt{2n}} \leq b\right) &= \sum_{x \in [a, b] \cap (2\mathbb{Z})/\sqrt{2n}} P\left(\frac{S_{2n}}{\sqrt{2n}} = x\right) \\ &\approx \sum_{x \in [a, b] \cap (2\mathbb{Z})/\sqrt{2n}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{\sqrt{2}}{\sqrt{n}}. \end{aligned}$$

The RHS approximates a Riemann sum of the function  $e^{-x^2/2}/\sqrt{2\pi}$  on  $[a, b]$  with grid size  $\sqrt{2}/\sqrt{n}$ . Since  $S_{2n+1} = S_{2n} \pm 1$ , the probability of  $S_{2n+1}/\sqrt{2n+1} \in [a, b]$  can be approximated by the same Riemann sum, thus we have proved

**Thm 98 (De Moivre-Laplace).** *If  $a < b$ ,  $m \rightarrow \infty$ , then*

$$P\left(a \leq \frac{S_m}{\sqrt{m}} \leq b\right) \rightarrow \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx.$$

**11.4. Weak convergence.** Recall that a sequence of distribution functions  $F_n$  converges weakly to a function  $F$  if  $F_n(x) \rightarrow F(x)$  for all

$x$  where  $F$  is continuous. A sequence of random variables  $X_n$  converges weakly to  $X$  if the associated distribution functions converge weakly.

It is often the case that weak convergence is defined in terms of probability measures. A sequence of probability measures  $\mu_n$  *converges weakly to a probability measure*  $\mu$  means the corresponding distribution function of  $\mu_n$  converges weakly to the distribution function of  $\mu$ .

**Example 43 (Exponential approximation to Geometric).**

*Suppose that a success trial happens with a probability proportional to time length and the probability is  $p$  ( $0 < p < 1$ ) per unit time. For  $t \geq 0$ , divide the interval  $[0, t]$  into pieces of equal size  $\delta > 0$ . Then the waiting time  $X$  for the first success follows  $\text{Geom}(\delta p)$ . Whence*

$$P(X > t) = (1 - \delta p)^{t/\delta} \rightarrow e^{-pt}, \text{ for } t \geq 0.$$

*So geometric random variable converges weakly to exponential. See Remark 8 for the interpretation of the parameter of exponential distribution.*



**Example 44 (Birthday problem).** Let  $X_1, X_2, \dots$  be i.i.d with uniform distribution on  $\{1, \dots, n\}$ , and

$$T_n = \min\{k : X_k = X_j \text{ for some } j < k\}$$

the first time some  $X_k$  repeats itself. Imagine  $X_k$ s are birthdays of people or coupons you collect (see Example 38), then  $T_n$  represents the first time people's birthday repeats or a type of coupon that you already have is collected. Clearly  $T_n$  takes values in  $\{2, \dots, n+1\}$ . For  $1 \leq k \leq n$ , the occurrence of the event  $T_n > k$  amounts to the first  $k$  observed values  $X_1, \dots, X_k$  being distinct, so

$$P(T_n > k) = \prod_{j=1}^k \left(1 - \frac{j-1}{n}\right).$$

It follows<sup>7</sup> that for  $x \geq 0$ ,  $k \approx n^{1/2}x$ , we have

$$P(T_n > n^{1/2}x) \rightarrow e^{-x^2/2} \text{ as } n \rightarrow \infty.$$

---

<sup>7</sup>Use Exercise 3.1.1 with  $c_{jn} = -\frac{j-1}{n}$  if  $j \leq k = \lceil n^{1/2}x \rceil$ , and 0 if  $j > k$ .

**Thm 99 (Scheffé Theorem).** *Suppose that  $p_n(x)$ ,  $p(x)$  are probability densities, and  $p_n \rightarrow p$ , a.s. Let*

$$\nu_n(A) = \int_A p_n(x)dx, \quad \nu(A) = \int_A p(x)dx$$

and

$$\|\nu_n - \nu\|_{TV} \triangleq \sup_{A \in \mathcal{B}(\mathbb{R})} |\nu_n(A) - \nu(A)|.$$

Then

$$\|\nu_n - \nu\|_{TV} \leq \int_{\mathbb{R}} |p_n - p| dx \rightarrow 0.$$

PROOF. The inequality follows from the property of integral. To see  $\int |p_n - p| \rightarrow 0$ , we note that, since

$$\int_{\mathbb{R}} (p - p_n) dx = 0,$$

we have

$$\int_{\mathbb{R}} (p - p_n)^+ dx = \int_{\mathbb{R}} (p - p_n)^- dx.$$

It follows that

$$\int_{\mathbb{R}} |p_n - p| dx = 2 \int_{\mathbb{R}} (p - p_n)^+ dx.$$

The RHS converges to zero, since  $(p - p_n)^+ \leq p$ ,  $p$  is integrable and  $(p - p_n)^+ \rightarrow 0$ , *a.s.* so that the dominated convergence theorem is applicable.  $\square$

REMARK 9. Given measures  $\nu_n, \nu$ ,  $\|\nu_n - \nu\|_{TV}$  defined in Theorem 99 is referred to as the **total variation norm** which gauges the discrepancy between  $\nu_n$  and  $\nu$ . Since

$$\sup_{A=(-\infty, x], x \in \mathbb{R}} |\nu_n(A) - \nu(A)| \leq \|\nu_n - \nu\|_{TV},$$

so convergence in total variation norm implies weak convergence of distributions. But the converse is not necessarily true. Consider the

example,

$$\nu_n = \delta_{1/n}, \quad \nu = \delta_0,$$

where  $\delta_a$  is the Dirac measure concentrated on  $a$ . Then  $\|\nu_n - \nu\|_{TV} = 1$  but for any  $x$ ,

$$\nu_n((-\infty, x]) = 1_{[1/n, \infty)}(x) \rightarrow 1_{[0, \infty)}(x) = \nu_0((-\infty, x]).$$

**Example 45 (Order statistics).** A sample of  $n$  points are picked randomly from the interval  $(0, 1)$ , i.e. the locations  $X_1, X_2, \dots, X_n$  are i.i.d with common uniform distribution on  $(0, 1)$ . The density of the  $k$ -th largest  $X_{(k)}$  is

$$p_k(x) = \frac{n!}{(k-1)!(n-k)!} x^{k-1} (1-x)^{n-k}.$$

PROOF. As in Example 41, we may assume that  $X_1, X_2, \dots, X_n$  are distinct. As a matter of fact, the probability of more than one points landing in a small interval of size  $\delta x$  is  $O((\delta x)^2)$ , which is already negligible. If  $X_{(k)} \in (x, x + \delta x)$ , then there are exactly  $k-1$  points whose location  $< x$ , and  $n-k$  points  $> x + \delta x$ . Imagine arranging  $n$

balls into the interval  $(0, 1)$ , one in the infinitesimal interval  $(x, x + \delta x)$  (there are  $n$  ways to pick one ball),  $k - 1$  to the left of it (there are  $C_{n-1}^{k-1}$  ways to pick  $k - 1$  balls out of  $n - 1$ ) and  $n - k$  to the right. The probability is

$$\begin{aligned} P(X_{(k)} \in (x, x + \delta x)) &= n C_{n-1}^{k-1} \cdot \delta x \cdot (1 - x - \delta x)^{n-k} \\ &= \frac{n!}{(k-1)!(n-k)!} x^{k-1} \cdot \delta x \cdot (1 - x - \delta x)^{n-k}. \end{aligned}$$

Therefore

$$p_k(x) = \lim_{\delta x \rightarrow 0} \frac{P(X_{(k)} \in (x, x + \delta x))}{\delta x} = \frac{n!}{(k-1)!(n-k)!} x^{k-1} (1-x)^{n-k}.$$

□

The preceding example shows that  $X_{(k)}$  has Beta distribution (Definition 46) with parameter  $k$  and  $n - k + 1$ , i.e.

$$p_k(x) = \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} x^{k-1} (1-x)^{n-k}, \quad 0 < x < 1.$$

Now assume that  $2n + 1$  points are picked randomly from  $(0, 1)$ . We will find the weak limit of the **central order statistics**  $X_{(n+1)}$  via Theorem 99. By Example 45,  $X_{(n+1)}$  has density

$$p_{n+1}(x) = \frac{(2n+1)!}{n!n!} x^n (1-x)^n \sim \text{Beta}(n+1, n+1).$$

Then by Example 36,

$$EX_{(n+1)} = \frac{1}{2}, \quad \text{Var}(X_{(n+1)}) = \frac{1}{4(2n+3)}.$$

Consider the standardization of  $X_{(n+1)}$ ,

$$Y_n = 2\sqrt{2n} \left( X_{(n+1)} - \frac{1}{2} \right).$$

Since we are interested in large  $n$  asymptotics, we have replaced  $(2n + 3)$  with  $2n$ . Through a change of variable  $x = 1/2 + y/(2\sqrt{2n})$ , the density of  $Y_n$  is found to be

$$\begin{aligned} p_{Y_n}(y) &= \frac{(2n+1)!}{n!n!} \left( \frac{1}{2} + \frac{y}{2\sqrt{2n}} \right)^n \left( \frac{1}{2} - \frac{y}{2\sqrt{2n}} \right)^n \frac{1}{2\sqrt{2n}} \\ &= C_{2n}^n 2^{-2n} \left( 1 - \frac{y^2}{2n} \right)^n \frac{2n+1}{2n} \frac{\sqrt{n}}{\sqrt{2}}. \end{aligned}$$

The factor  $C_{2n}^n 2^{-2n}$  is identified to be  $P(S_{2n} = 0)$  from Lemma 97, thus asymptotic to  $1/\sqrt{\pi n}$ , hence

$$p_{Y_n}(y) \rightarrow \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}}.$$

It follows from Theorem 99 that  $Y_n$  converges weakly to the normal distribution  $\mathcal{N}(0, 1)$ .

LEMMA 100. *Let  $g$  be a measurable function from  $\mathbb{R}^n$  to  $\mathbb{R}$ ,  $D$  the set of discontinuous points of  $g$ . Then  $D$  is Borel measurable, i.e.  $D \in \mathcal{B}(\mathbb{R}^n)$ .*

PROOF. For any positive rationals  $\varepsilon, \rho$ , let

$$D(\varepsilon, \rho) = \{x \in \mathbb{R}^n : \exists y, z \in B_\rho(x) \text{ such that } |g(y) - g(z)| \geq \varepsilon\},$$

where  $B_\rho(x) = \{x' : |x' - x| < \rho\}$ . Then

$$D = \bigcup_{\varepsilon} \bigcap_{\rho} D(\varepsilon, \rho).$$

We claim that  $D(\varepsilon, \rho)$  is open so that the conclusion follows immediately. Indeed, let  $x \in D(\varepsilon, \rho)$ , then there are  $y, z \in B_\rho(x)$  such that  $|g(y) - g(z)| \geq \varepsilon$ . There is an open ball  $B_{\rho_1}(x)$  such that every point in it has a distance less than  $\rho$  from  $y$ , and there is another open ball  $B_{\rho_2}(x)$  such that every point in it has a distance less than  $\rho$  from  $z$ . Now take the intersection  $B_{\rho_1 \wedge \rho_2}(x)$  of the two open balls. Clearly for every  $x'$  in  $B_{\rho_1 \wedge \rho_2}(x)$ , we have that  $y, z \in B_\rho(x')$  and  $|g(y) - g(z)| \geq \varepsilon$ .



Hence

$$B_{\rho_1 \wedge \rho_2}(x) \subset D(\varepsilon, \rho).$$

This confirms that  $D(\varepsilon, \rho)$  is an open set. □

**Thm 101 (Continuous mapping theorem).** *Let  $g$  be a measurable function from  $\mathbb{R}$  to  $\mathbb{R}$ ,  $D_g$  the set of discontinuities of  $g$ . If  $X_n \rightarrow X$  weakly and  $P(X \in D_g) = 0$ , then*

$$g(X_n) \rightarrow g(X) \text{ weakly.}$$

*If additionally  $g$  is bounded, then  $Eg(X_n) \rightarrow Eg(X)$ .*

REMARK 10. By Lemma 100, it makes sense to write  $P(X \in D_g)$ .

PROOF. 1. In view of Skorohod theorem 26, there are  $Y_n, Y$  defined on a common probability space  $(\Omega, \mathcal{F}, \mu)$  so that

$$Y_n \stackrel{d}{=} X_n, Y \stackrel{d}{=} X \text{ and } Y_n \rightarrow Y, a.s.$$

Since  $Y$  and  $X$  have identical distribution, we get that

$$\mu(Y \in B) = P(X \in B) \text{ for any } B \in \mathcal{B}(\mathbb{R}).$$

Setting  $B = D_g$  gives  $\mu(Y \in D_g) = 0$ . Thus the almost sure convergence  $Y_n \rightarrow Y$ , *a.s.* remains true outside  $D_g$ . But  $g$  is continuous in the complement of  $D_g$ , so overall we have  $g(Y_n) \rightarrow g(Y)$  *a.s.* which implies that  $g(Y_n) \rightarrow g(Y)$  weakly (Theorem 24). If  $z$  is a continuity point of  $g(X)$ , i.e.  $P(g(X) = z) = 0$ , then using that  $g(Y) \stackrel{d}{=} g(X)$ , we see that  $z$  is also a continuity point of  $g(Y)$ . Hence

$$\mu(g(Y_n) \leq z) \rightarrow \mu(g(Y) \leq z) = P(g(X) \leq z).$$

Now using that  $g(Y_n) \stackrel{d}{=} g(X_n)$ , we have

$$\mu(g(Y_n) \leq z) = P(g(X_n) \leq z).$$

It follows that

$$P(g(X_n) \leq z) \rightarrow P(g(X) \leq z).$$

Therefore  $g(X_n) \rightarrow g(X)$  weakly.

**2.** If  $g$  is bounded, then bounded convergence theorem gives  $Eg(Y_n) \rightarrow Eg(Y)$ . By the identical distribution construction, we have  $Eg(Y_n) = Eg(X_n)$  and  $Eg(Y) = Eg(X)$ , thus  $Eg(X_n) \rightarrow Eg(X)$ .  $\square$

Weak convergence can be characterized through dual actions with bounded continuous functions.

**Thm 102.** *Let  $X_n$  and  $X$  be random variables. The following are equivalent.*

- (1)  $X_n \rightarrow X$  weakly.
- (2) For every bounded continuous function  $g$ , it holds that

$$Eg(X_n) \rightarrow Eg(X).$$

- (3) For every ***P*-continuity set**  $A$ , i.e.

$$A \in \mathcal{B}(\mathbb{R}) \text{ and } P(X \in \partial A) = 0,$$

*it holds that*

$$P(X_n \in A) \rightarrow P(X \in A).$$

**PROOF. 1.** (1)  $\Rightarrow$  (2) follow from the continuous mapping theorem (Theorem 101).

**2.** To see (1)  $\Rightarrow$  (3), continue the proof of Theorem 101 and take  $g = 1_A$ , then  $D_g = \partial A$ . If  $P(X \in \partial A) = 0$ , then the random

variable  $Y$  constructed there has  $\mu(Y \in D_g) = 0$ , hence we still have  $g(Y_n) \rightarrow g(Y)$  *a.s.* It follows that

$$\mu(Y_n \in A) = Eg(Y_n) \rightarrow Eg(Y) = \mu(Y \in A).$$

By the identical distribution construction, this translates to

$$P(X_n \in A) \rightarrow P(X \in A).$$

**3.** (3)  $\Rightarrow$  (1). For every  $A = (-\infty, x]$ ,  $\partial A = \{x\}$ ,  $P(X \in \partial A) = 0$  means that  $x$  is a continuous point of the distribution function of  $X$ , thus  $P(X_n \in A) \rightarrow P(X \in A)$  is the same thing as the distribution function of  $X_n$  being convergent to that of  $X$  at  $x$ .

**4.** (2)  $\Rightarrow$  (1). Fix  $x \in \mathbb{R}$ . For  $y > x$ , consider the continuous piecewise linear function  $g(s)$  which equal  $1_{(-\infty, x]}$  if  $s \leq x$ , linear on  $s \in [x, y]$  and equals 0 on  $s \geq y$ . Then

$$P(X_n \leq x) = E1_{(-\infty, x]}(X_n) \leq Eg(X_n).$$

Hence

$$\limsup_n P(X_n \leq x) \leq \limsup_n Eg(X_n) = Eg(X) \leq P(X \leq y).$$

Now let  $y \downarrow x$  gives

$$\limsup_n P(X_n \leq x) \leq P(X \leq x).$$

Similarly, for  $z < x$ , consider the continuous piecewise linear function  $h(s)$  which equal  $1_{(-\infty, z]}$  if  $s \leq z$ , linear on  $s \in [z, x]$  and equals 0 on  $s \geq x$ . Then

$$P(X_n \leq x) \geq Eh(X_n).$$

Hence

$$\liminf_n P(X_n \leq x) \geq \liminf_n Eh(X_n) = Eh(X) \geq P(X \leq z).$$

Now let  $z \uparrow x$  gives

$$\liminf_n P(X_n \leq x) \geq P(X < x).$$

Putting these inequalities together, we have

$$P(X < x) \leq \liminf_n P(X_n \leq x) \leq \limsup_n P(X_n \leq x) \leq P(X \leq x).$$

So the weak convergence follows. □

**Thm 103 (Helly's selection theorem).** *For every sequence of distribution functions  $F_n$ , there exist a subsequence  $F_{n_k}$  and a non-decreasing, right-continuous function  $F$  such that  $F_{n_k} \rightarrow F$  at every continuity point of  $F$ .*

PROOF. Note  $0 \leq F_n \leq 1$ ,  $\forall n$ . Following a standard diagonal procedure, we may find a subsequence  $F_{n_k}$  which converges at every rational number  $q \in \mathbb{Q}$  to some  $G(q)$ . Clearly  $G$  is nondecreasing in  $\mathbb{Q}$ . Let

$$F(x) = \inf\{G(q) : x < q\}, \quad \forall x.$$

Then  $F$  is nondecreasing.

1.  $F$  is right continuous at any  $x$ . For  $\varepsilon > 0$ , there is  $q \in \mathbb{Q}$ ,  $x < q$  such that

$$F(x) \leq G(q) < F(x) + \varepsilon.$$

If  $x \leq y < q$ , then  $F(y) \leq G(q)$ . It follows that

$$F(y) \leq G(q) < F(x) + \varepsilon.$$

Hence  $F$  is right continuous.

**2.** If  $F$  is continuous at  $x$ , then  $F_{n_k}(x) \rightarrow F(x)$  hence the conclusion follows. To see this, first note that by continuity, for  $\varepsilon > 0$  there is  $y < x$  such that

$$F(x) - \varepsilon < F(y).$$

Now choose rational numbers  $q, r$  so that  $y < q < x < r$  and  $G(r) < F(x) + \varepsilon$ . Then

$$F(x) - \varepsilon < F(y) \leq G(q) \leq G(r) < F(x) + \varepsilon.$$

By the definition of  $G$ ,

$$F_{n_k}(q) \rightarrow G(q), \quad F_{n_k}(r) \rightarrow G(r).$$

Therefore using the monotonicity of  $F_{n_k}$ , for large  $k$ ,

$$F(x) - \varepsilon < F_{n_k}(q) \leq F_{n_k}(x) \leq F_{n_k}(r) < F(x) + \varepsilon.$$

It follows that

$$F(x) - \varepsilon \leq \liminf_k F_{n_k}(x) \leq \limsup_k F_{n_k}(x) \leq F(x) + \varepsilon.$$

Finally let  $\varepsilon \rightarrow 0$ . □

REMARK 11. The limit function  $F$  of Helly's selection theorem necessarily has  $0 \leq F \leq 1$ , but the theorem does not claim that  $F$  is a distribution function, the reason being that probability mass could escape to infinity. For example,  $F_n$  is the distribution function of a uniform random variable on  $(n, n+1)$ , clearly  $F_n \rightarrow F \equiv 0$ , but  $F$  is not a distribution function. A necessary and sufficient condition to avoid this situation is **tightness**.

**Def 47.** A sequence of probability measures  $\mu_n$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  is **tight** if  $\forall \varepsilon$ , there is an interval  $(a, b]$  such that

$$\mu_n((a, b]) > 1 - \varepsilon, \quad \forall n.$$

In terms of the corresponding distribution functions  $F_n$ , this means that  $\forall \varepsilon$ , there is  $M_\varepsilon > 0$  such that

$$\limsup_n F_n(-M_\varepsilon) + 1 - F_n(M_\varepsilon) \leq \varepsilon.$$

**Thm 104.** A sequence of probability measures  $\mu_n$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  is tight if and only if every subsequence has a further subsequence that is convergent weakly to some probability measure.



**PROOF. 1.** Suppose that  $\{\mu_n\}$  is tight. Denote by  $F_n$  the distribution function corresponding to  $\mu_n$ . In view of Helly's selection theorem, there exist a subsequence  $F_{n_k}$  and a nondecreasing, right-continuous function  $F$  such that  $F_{n_k} \rightarrow F$  at every continuity point of  $F$ . Then there is a unique measure  $\mu$  such that  $\mu((a, b]) = F(b) - F(a)$ ,  $\forall a, b$ . It remain to prove that  $\mu$  is a probability measure. Due to tightness,  $\forall \varepsilon > 0$ , we can find  $(a, b]$  so that

$$\mu_{n_k}((a, b]) > 1 - \varepsilon, \quad \forall k.$$

Additionally we may increase  $b$  and decrease  $a$  so that they are continuity points of  $F$ . Now taking limit in the equation

$$F_{n_k}(b) - F_{n_k}(a) = \mu_{n_k}((a, b]) > 1 - \varepsilon$$

yields  $\mu((a, b]) \geq 1 - \varepsilon$ . But  $\varepsilon$  is arbitrary, hence  $\mu$  is a probability measure.

**2.** Suppose the stated subsequential property holds but  $\{\mu_n\}$  is not tight. Then there is  $\varepsilon_0 > 0$  so that for every interval  $(a, b]$ , there

is an index  $n'$  having

$$\mu_{n'}((a, b]) \leq 1 - \varepsilon_0.$$

This is particularly true for the sequence of intervals  $(-k, k]$ , i.e. there is a subsequence  $n_k$  so that

$$\mu_{n_k}((-k, k]) \leq 1 - \varepsilon_0, \quad \forall k.$$

According to the stated subsequential property, there is a subsequence of  $\{n_k\}$ , say  $\{n_{k_l}\}$ , such that  $\mu_{n_{k_l}}$  converges weakly to some probability measure  $\mu$ . Choose  $(a, b]$  with  $\mu((a, b]) > 1 - \varepsilon_0$  and  $a, b$  being the continuity point of  $\mu$ , i.e.  $\mu(\{a\}) = \mu(\{b\}) = 0$ . Hence  $\mu_{n_{k_l}}((a, b]) \rightarrow \mu((a, b])$  as  $l \rightarrow \infty$ . Therefore as soon as  $l$  is large, we can have the interval  $(a, b]$  contained in  $(-k_l, k_l]$  and

$$1 - \varepsilon_0 < \mu_{n_{k_l}}(a, b] \leq \mu_{n_{k_l}}((-k_l, k_l]) \leq 1 - \varepsilon_0.$$

A contradiction. □

**COROLLARY 5.** *Let  $\mu_n$  be a tight sequence of probability measures  $\mu_n$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ ,  $\mu$  a probability measure. If every subsequence of*

$\mu_n$  that converges weakly actually converges to  $\mu$ , then  $\mu_n$  converges weakly to  $\mu$ .

PROOF. By Theorem 104 and the assumption, every subsequence has a further subsequence that converges weakly to  $\mu$ . Suppose that  $\mu_n$  does not converge weakly to  $\mu$ . Then there is  $x$  with  $\mu(\{x\}) = 0$  but  $\mu_n((-\infty, x])$  does not converge to  $\mu((-\infty, x])$ , i.e., there exists  $\varepsilon_0 > 0$  so that along a subsequence  $\mu_{n_k}$ ,

$$|\mu_{n_k}((-\infty, x]) - \mu((-\infty, x])| \geq \varepsilon_0.$$

Then no subsequence of  $\mu_{n_k}$  can converge weakly to  $\mu$ , a contradiction to the assumption.  $\square$

Below is a sufficient condition for tightness.

**Thm 105.** Suppose  $\varphi \geq 0$  and  $\lim_{|x| \rightarrow \infty} \varphi = \infty$ . If

$$C = \sup_n \int \varphi(x) dF_n(x) < \infty,$$

then  $F_n$  is tight.

PROOF. Write  $\mu_n$  for the probability measure corresponding to  $F_n$  (Recall Remark 3). Consider the interval  $(a, b]$ . We have

$$\begin{aligned}\mu_n((a, b]^c) &= \int 1_{(a, b]^c}(x) d\mu_n(x) = \int_{(a, b]^c} \frac{1}{\varphi(x)} \varphi(x) d\mu_n(x) \\ &\leq \frac{1}{\inf_{x \in (a, b]^c} \varphi} \int \varphi(x) d\mu_n(x) \leq \frac{C}{\inf_{x \in (a, b]^c} \varphi}.\end{aligned}$$

By the assumption, the rightmost term converges to zero as  $a \rightarrow -\infty$ ,  $b \rightarrow \infty$ . Therefore for  $\varepsilon > 0$ , there are  $a, b$  so that

$$\mu_n((a, b]) > 1 - \varepsilon, \quad \forall n.$$

□

**Example 46.** Let  $k \geq 1$  and  $X_n$  have distribution function  $F_n$ . If  $C = \sup_n E|X_n|^k < \infty$ , then  $F_n$  is tight.

## 12. Characteristic function

### 12.1. A brief review of complex calculus.

**Def 48.** For the complex valued random variable,

$$X = Y + iZ,$$

its expectation is defined to be

$$EX = EY + iEZ.$$

As in the real case, the complex exponential can be defined via Taylor series.

**Def 49.** The complex exponential is defined to be

$$e^z = \sum_{k=0}^{\infty} \frac{z^k}{k!}, \quad z \in \mathbb{C}.$$

If  $z = ix$  with  $x \in \mathbb{R}$  being a real number, then the complex exponential becomes the expansion

$$e^{ix} = \sum_{k=0}^{\infty} \frac{(-1)^k}{(2k)!} x^{2k} + i \sum_{k=0}^{\infty} \frac{(-1)^k}{(2k+1)!} x^{2k+1},$$

which gives the frequently used **Euler's formula**:

$$e^{ix} = \cos x + i \sin x.$$

## 12.2. The definition of Characteristic function.

**Def 50.** *The characteristic function of a probability measure  $\mu$  on  $\mathbb{R}$  is defined to be*

$$\varphi(t) = \int_{\mathbb{R}} e^{itx} d\mu(x) = \int_{\mathbb{R}} \cos(tx) d\mu(x) + i \int_{\mathbb{R}} \sin(tx) d\mu(x).$$

*The characteristic function of a random variable  $X$  with distribution  $\mu$  is*

$$\varphi(t) = Ee^{itX} = \int_{\mathbb{R}} e^{itx} d\mu(x).$$

Note the complex number  $e^{itX}$  has a bounded modulus:  $|e^{itX}| = 1$ , hence the characteristic function of a probability measure or random variable always exists.

**Thm 106.** *All characteristic functions have the following properties.*

$$(1) \quad \varphi(0) = 1; \quad \varphi(-t) = \overline{\varphi(t)};$$

(2)

$$|\varphi(t)| = |Ee^{itX}| \leq E|e^{itX}| = 1;$$

(3)

$$Ee^{it(aX+b)} = e^{itb}\varphi(at);$$

(4)  $\varphi(t)$  is uniformly continuous on  $\mathbb{R}$ , moreover

$$|\varphi(t+h) - \varphi(t)| \leq E|e^{ihX} - 1|.$$

PROOF. Omitted. □

**Thm 107 (Convolution).** *If  $X$  and  $Y$  are independent with respective characteristic functions  $\varphi_X, \varphi_Y$ , then  $X + Y$  has the characteristic function  $\varphi_{X+Y} = \varphi_X \cdot \varphi_Y$ .*

PROOF.

$$Ee^{it(X+Y)} = E(e^{itX} \cdot e^{itY}) = Ee^{itX} \cdot Ee^{itY}.$$

□

**Example 47 (Ch. of Normal distribution).**  $X \sim \mathcal{N}(0, 1)$ ,

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad \varphi(t) = e^{-\frac{t^2}{2}}.$$

Hence  $\sigma X + \mu \sim \mathcal{N}(\mu, \sigma^2)$  has

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad \varphi(t) = e^{i\mu t - \frac{\sigma^2 t^2}{2}}.$$

PROOF. To avoid working in the complex domain, we adopt an indirect method that has the advantage of being elementary. By the definition and Euler formula,

$$\begin{aligned} \varphi(t) &= \frac{1}{\sqrt{2\pi}} \int e^{itx} e^{-\frac{1}{2}x^2} dx = \frac{1}{\sqrt{2\pi}} \int e^{-\frac{1}{2}x^2} (\cos tx + i \sin tx) dx \\ &= \frac{1}{\sqrt{2\pi}} \int e^{-\frac{1}{2}x^2} \cos tx dx. \end{aligned}$$



By dominated convergence theorem and integration by parts,

$$\begin{aligned}\varphi'(t) &= -\frac{1}{\sqrt{2\pi}} \int e^{-\frac{1}{2}x^2} x \sin tx dx = \frac{1}{\sqrt{2\pi}} \int \sin tx d\left(e^{-\frac{1}{2}x^2}\right) \\ &= -\frac{1}{\sqrt{2\pi}} \int t e^{-\frac{1}{2}x^2} \cos tx dx = -t\varphi(t).\end{aligned}$$

The equation

$$\varphi'(t) = -t\varphi(t) \text{ with } \varphi(0) = 1$$

has a unique solution

$$\varphi(t) = e^{-\frac{t^2}{2}}.$$

The characteristics function of  $\sigma X + \mu$  is obtained from Theorem 106.

□

**12.3. The inversion formula and uniqueness.** The inversion formula below indicates that probability measures with identical characterisitic function must be equal. But before proceeding to prove the

formula, it is helpful to recall the fact from analysis,

$$\int_0^{\infty} \frac{\sin x}{x} dx = \frac{\pi}{2}.$$

A variant of this formula will be used in the proof of the inversion formula. For  $a \in \mathbb{R}$ ,

$$(12.1) \quad \int_{-\infty}^{\infty} \frac{\sin ax}{x} dx = 2 \int_0^{\infty} \frac{\sin ax}{x} dx = \text{sign}(a) \cdot \pi,$$

where

$$\text{sign}(a) = \begin{cases} -1, & \text{if } a < 0, \\ 0, & \text{if } a = 0, \\ 1, & \text{if } a > 0. \end{cases}$$

**Thm 108 (The inversion formula).** *Let  $\varphi$  be the characteristic function of the probability measure  $\mu$ . Then for  $a < b$ ,*

$$\mu((a, b)) + \frac{1}{2}\mu(\{a, b\}) = \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt.$$

PROOF. Let

$$I_T = \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{2\pi it} \varphi(t) dt = \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{2\pi it} \left( \int_{\mathbb{R}} e^{itx} d\mu(x) \right) dt.$$

Since the integrand is bounded,

$$\left| \frac{e^{-ita} - e^{-itb}}{it} e^{itx} \right| \leq \left| \frac{e^{-ita} - e^{-itb}}{it} \right| = \left| \int_a^b e^{-itx} dx \right| \leq |b - a|,$$

we can employ Fubini theorem to write

$$I_T = \int_{\mathbb{R}} \left( \int_{-T}^T \frac{e^{it(x-a)} - e^{it(x-b)}}{2\pi it} dt \right) d\mu(x).$$

Since

$$\begin{aligned} & e^{it(x-a)} - e^{it(x-b)} \\ &= \cos(t(x-a)) + i \sin(t(x-a)) - \cos(t(x-b)) - i \sin(t(x-b)), \end{aligned}$$

the real part is even in  $t$ . It follows that

$$I_T = \int_{\mathbb{R}} \left( \int_{-T}^T \frac{\sin(t(x-a))}{2\pi t} dt - \int_{-T}^T \frac{\sin(t(x-b))}{2\pi t} dt \right) d\mu(x).$$

From Equation (12.1), we see that, as soon as  $T$  is larger than some  $T_0 > 0$ , the bracketed integrand is bounded by

$$\begin{aligned} & \left| \int_{-T}^T \frac{\sin(t(x-a))}{2\pi t} dt \right| + \left| \int_{-T}^T \frac{\sin(t(x-b))}{2\pi t} dt \right| \\ & \leq \frac{1}{2} |\text{sign}(x-a)| + \frac{1}{2} |\text{sign}(x-b)| + 1 \leq 2. \end{aligned}$$

Hence we can invoke the bounded convergence theorem to conclude that

$$\lim_{T \rightarrow \infty} I_T = \int_{\mathbb{R}} \left( \frac{1}{2} \text{sign}(x-a) - \frac{1}{2} \text{sign}(x-b) \right) d\mu(x).$$

Now

$$\frac{1}{2}\text{sign}(x-a) - \frac{1}{2}\text{sign}(x-b) = \begin{cases} 0, & \text{if } x < a \text{ or } x > b, \\ 1/2, & \text{if } x = a \text{ or } x = b, \\ 1, & \text{if } a < x < b. \end{cases}$$

Therefore

$$\lim_{T \rightarrow \infty} I_T = \mu((a, b)) + \frac{1}{2}\mu(\{a, b\}).$$

□

**Thm 109 (Uniqueness).** *The inversion formula implies uniqueness, i.e. if  $\mu, \nu$  have identical characteristic function, then  $\mu$  and  $\nu$  must coincide.*

**PROOF.** If  $\mu, \nu$  have identical characteristic function  $\varphi$ , then by the inversion formula

$$\mu((a, b]) = \nu((a, b]) = \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt, \quad \forall a < b,$$

provided  $\mu(\{a, b\}) = \nu(\{a, b\}) = 0$ . So  $\mu$  and  $\nu$  coincide on the class  $\mathcal{S}_0$  of intervals whose extreme points do not hold mass,

$$\mathcal{S}_0 = \{(a, b) : a, b \in \mathbb{R}, a, b \notin A\},$$

where

$$A = \{a \in \mathbb{R} : \mu(\{a\}) \text{ or } \nu(\{a\}) \neq 0\}.$$

But  $A$  is at most countable, hence  $A^c$  is dense. It follows that  $\mu$  and  $\nu$  coincide on the  $\pi$ -system  $\{(a, b) : a, b \in \mathbb{R}\}$ . Hence  $\mu$  equals  $\nu$  by uniqueness (Theorem 10).  $\square$

**COROLLARY 6.** *If  $\varphi_X$  is real, then  $X$  and  $-X$  have the same distribution.*

**COROLLARY 7 (Sum of normal distributions).** *Let*

$$X \sim \mathcal{N}(\mu_1, \sigma_1^2), Y \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

*be independent, then*

$$\varphi_{X+Y} = e^{i\mu_1 t - \sigma_1^2 t^2 / 2} \cdot e^{i\mu_2 t - \sigma_2^2 t^2 / 2} = \exp \left[ i(\mu_1 + \mu_2)t - \frac{(\sigma_1^2 + \sigma_2^2)t^2}{2} \right],$$

which is recognized as the characteristic function of  $\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$ . Combining this with the uniqueness, we see that

$$X + Y \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2).$$

Compare Example 34.

**Thm 110 (Integrable Ch.f.).** *If  $\mu$  has integrable characteristic function  $\varphi$ , then it admits bounded continuous density,*

$$p(x) = \frac{1}{2\pi} \int e^{-itx} \varphi(t) dt.$$

PROOF. As in the proof of the inversion formula,

$$\left| \frac{e^{-ita} - e^{-itb}}{it} e^{itx} \right| \leq |b - a|.$$

Together with the assumption that  $\varphi$  is integrable, this says

$$\frac{e^{-ita} - e^{-itb}}{it} \varphi(t)$$

is integrable, therefore by the dominated convergence theorem, the inversion formula takes the form,  $\forall a < b$ ,

$$\begin{aligned}\mu((a, b)) + \frac{1}{2}\mu(\{a, b\}) &= \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt,\end{aligned}$$

which gives

$$\frac{1}{2}\mu(\{a\}) \leq \mu((a, b)) + \frac{1}{2}\mu(\{a, b\}) \leq \frac{|b - a|}{2\pi} \int_{-\infty}^{\infty} |\varphi(t)| dt.$$



Letting  $b \rightarrow a$ , this shows that  $\mu(\{a\}) = 0$ , i.e.,  $\mu$  does not have point mass, and it follows that by Fubini theorem,  $\forall h > 0$ ,

$$\begin{aligned}\mu((a, a+h)) &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-ita} - e^{-it(a+h)}}{it} \varphi(t) dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \left( \int_a^{a+h} e^{-itx} dx \right) \varphi(t) dt \\ &= \int_a^{a+h} \left( \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt \right) dx.\end{aligned}$$

Hence  $\mu$  has density

$$p(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt.$$

Using again that  $\varphi$  is integrable and the dominated convergence theorem,  $p(x)$  is continuous in  $x$ .  $\square$

The theorem indicates the mutually inversive relation between probability measure and characteristic function: if  $\mu$  has density  $p(x)$

and its characterisitic function  $\varphi(t)$  is integrable, then

$$\varphi(t) = \int e^{itx} p(x) dx \iff p(x) = \frac{1}{2\pi} \int e^{-itx} \varphi(t) dt.$$

## 12.4. Moments and derivatives.

LEMMA 111. *Let  $n \geq 0$ . For any  $x \in \mathbb{R}$ ,*

$$\left| e^{ix} - \sum_{k=0}^n \frac{(ix)^k}{k!} \right| \leq \min \left\{ \frac{2|x|^n}{n!}, \frac{|x|^{n+1}}{(n+1)!} \right\}.$$

PROOF. Integrating by parts gives, for  $n \geq 0$ ,

$$\begin{aligned} \int_0^x (x-s)^n e^{is} ds &= -\frac{1}{n+1} \int_0^x e^{is} d(x-s)^{n+1} \\ &= \frac{x^{n+1}}{n+1} + \frac{i}{n+1} \int_0^x (x-s)^{n+1} e^{is} ds. \end{aligned}$$

Since  $e^{ix} = 1 + i \int_0^x e^{is} ds$ , we get by induction that

$$e^{ix} = \sum_{k=0}^n \frac{(ix)^k}{k!} + \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds.$$

Hence the expansion error term has

$$\left| \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds \right| \leq \frac{i^{n+1}}{n!} \left| \int_0^x |(x-s)^n e^{is}| ds \right| \leq \frac{|x|^{n+1}}{(n+1)!}.$$

Integrating by parts again, the error term can also be written as,

$$\begin{aligned} \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds &= \frac{i^n}{n!} \int_0^x (x-s)^n de^{is} \\ &= -\frac{i^n x^n}{n!} + \frac{i^n}{(n-1)!} \int_0^x (x-s)^{n-1} e^{is} ds \\ &= \frac{i^n}{(n-1)!} \int_0^x (x-s)^{n-1} (e^{is} - 1) ds. \end{aligned}$$

Noticing  $|e^{is} - 1| \leq 2$ ,

$$\left| \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds \right| \leq \frac{2|x|^n}{n!}.$$

Putting together both estimates of the error term, we obtain

$$\left| \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds \right| \leq \min \left\{ \frac{2|x|^n}{n!}, \frac{|x|^{n+1}}{(n+1)!} \right\}.$$

The first term on the RHS is good for large  $|x|$ , while the second for small  $|x|$ .  $\square$

An immediate consequence of Lemma 111 is the error estimate of the Taylor expansion of characteristic function.

**Thm 112 (Taylor formula error estimate).** *Let  $n \geq 0$ . If  $E|X|^n < \infty$ , then*

$$\left| \varphi(t) - \sum_{k=0}^n \frac{(it)^k}{k!} EX^k \right| \leq E \left( \min \left\{ \frac{2|tX|^n}{n!}, \frac{|tX|^{n+1}}{(n+1)!} \right\} \right).$$

Despite the appearance of  $|X|^{n+1}$  on the RHS, the error term is indeed bounded from above by a constant multiple of  $E|tX|^n$  (without any requirement on the next order moment  $E|tX|^{n+1}$ ), since we always have  $\min\{\cdot, \cdot\} \leq 2|tX|^n/n! \leq |tX|^n$ . This simple observation readily leads us to the differentiability of  $\varphi$ .

**Thm 113.** *If  $E|X|^2 < \infty$ , then*

$$\varphi(t) = 1 + itEX + \frac{(it)^2}{2}EX^2 + o(t^2).$$

*This says  $\varphi$  is twice differentiable at  $t = 0$ ,*

$$\varphi'(0) = iEX, \quad \varphi''(0) = i^2EX^2.$$

More generally, existence of moment implies differentiability.

**Thm 114.** *Let  $k \geq 1$ . If  $\mu$  has  $\int |x|^k d\mu(x) < \infty$ , then its characteristic function  $\varphi$  has continuous derivative of order  $k$ ,*

$$\varphi^{(k)}(t) = \int (ix)^k e^{itx} d\mu(x).$$

Conversely, under mild conditions, existence of derivative at  $t = 0$  implies existence of moment. The following is a statement for second derivative, the case for higher even order derivative can be proved by induction.

**Thm 115.** *If*

$$\limsup_{h \rightarrow 0+} \frac{\varphi(h) - 2\varphi(0) + \varphi(-h)}{h^2} > -\infty,$$

*then*  $E|X|^2 < \infty$ .

PROOF. For  $h > 0$ , writing down the second order central finite difference of  $t \mapsto e^{itx}$  at the origin gives

$$\frac{e^{ihx} - 2 + e^{-ihx}}{h^2} = -2 \frac{1 - \cos(hx)}{h^2} \leq 0.$$

Denote by  $F$  the distribution function of  $X$ . By Fatou lemma

$$\begin{aligned}\int x^2 dF(x) &= \int \liminf_{h \rightarrow 0+} \left( 2 \frac{1 - \cos(hx)}{h^2} \right) dF(x) \\ &\leq \liminf_{h \rightarrow 0+} \int 2 \frac{1 - \cos(hx)}{h^2} dF(x) \\ &= - \limsup_{h \rightarrow 0+} \int \frac{e^{ihx} - 2 + e^{-ihx}}{h^2} dF(x) < \infty.\end{aligned}$$

□

Next is a simple condition which ensures that the Taylor expansion holds.

LEMMA 116. *If  $t \in \mathbb{R}$  satisfies*

$$\lim_{n \rightarrow \infty} \frac{|t|^n E|X|^n}{n!} = 0,$$

then  $\varphi(t)$  must have the expansion,

$$\varphi(t) = \sum_{k=0}^{\infty} \frac{(it)^k}{k!} EX^k.$$

This is particularly true if  $t$  has

$$\sum_{k=0}^{\infty} \frac{|t|^k}{k!} EX^k = Ee^{|tX|} < \infty.$$

PROOF. By Theorem 112, for each  $n$ , the error is no more than  $2|t|^n E|X|^n/n!$ . Hence under the assumption, the error goes to zero as  $n \rightarrow \infty$ .  $\square$

We can derive the characteristic functions of some distributions that are normally calculated via complex contour integration.

**Example 48 (Ch.f. of Normal distribution).** *Note in the first place that the normal distribution  $\mathcal{N}(0, 1)$  satisfies  $Ee^{|tX|} < \infty, \forall t$ .*



Therefore

$$\varphi(t) = \sum_{k=0}^{\infty} \frac{(it)^k}{k!} EX^k = \sum_{k=0}^{\infty} \frac{(it)^{2k}}{(2k)!} EX^{2k} = \sum_{k=0}^{\infty} \frac{(-t^2)^k}{(2k)!!} = e^{-\frac{t^2}{2}}.$$

**Example 49 (Ch.f. of Gamma distribution).**  $X \sim \text{Gamma}(\alpha, 1)$ ,  $\alpha > 0$ . Then,

$$\begin{aligned} \varphi(t) &= \frac{1}{\Gamma(\alpha)} \sum_{k=0}^{\infty} \frac{(it)^k}{k!} \int x^{k+\alpha-1} e^{-x} dx = \sum_{k=0}^{\infty} \frac{\Gamma(\alpha+k)}{\Gamma(\alpha)} \frac{(it)^k}{k!} \\ &= \sum_{k=0}^{\infty} \binom{-\alpha}{k} \frac{(-it)^k}{k!} = (1-it)^{-\alpha}. \end{aligned}$$

## 12.5. Continuity theorem.

**Thm 117 (Continuity theorem).** Let  $\mu_n, \mu$  be probability measures with characterisitic functions  $\varphi_n, \varphi$ . Then

$$\mu_n \rightarrow \mu \text{ weakly} \iff \varphi_n(t) \rightarrow \varphi(t), \forall t.$$

PROOF. 1. " $\Rightarrow$ ". By the definition

$$\varphi_n(t) = \int \cos(tx) d\mu_n(x) + i \int \sin(tx) d\mu_n(x).$$

Since  $\cos(tx)$ ,  $\sin(tx)$  are bounded continuous in  $x$ , Theorem 102 applies respectively to the real and imaginary part, hence

$$\varphi_n(t) \rightarrow \int \cos(tx) d\mu(x) + i \int \sin(tx) d\mu(x) = \varphi(t).$$

2. " $\Leftarrow$ ". If  $\mu_n$  is tight, then by Theorem 104 every subsequence has a further subsequence  $\mu_{n_k}$  which converges weakly to some probability measure  $\nu$ . By the already proved " $\Rightarrow$ " part, the characteristic function of  $\nu$  must equal  $\lim_k \varphi_{n_k}$ . Hence by Corollary 5 and Theorem 109,  $\nu$  coincides  $\mu$ . It then follows that  $\mu_n \rightarrow \mu$  weakly. Therefore it suffices

to prove the tightness of  $\mu_n$ . For this, by Fubini theorem,

$$\begin{aligned}\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi_n(t)) dt &= \int_{-\infty}^{\infty} \left( \frac{1}{\delta} \int_{-\delta}^{\delta} (1 - e^{itx}) dt \right) d\mu_n(x) \\ &= 2 \int_{-\infty}^{\infty} \left( 1 - \frac{\sin(\delta x)}{\delta x} \right) d\mu_n \\ &\geq 2 \int_{|x| \geq 2/\delta} \left( 1 - \frac{1}{|\delta x|} \right) d\mu_n \geq \mu_n(|x| \geq 2/\delta).\end{aligned}$$

Since  $\varphi$  is continuous at  $t = 0$  and  $\varphi(0) = 1$ , for  $\varepsilon > 0$  there is  $\delta > 0$ ,

$$\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi(t)) dt \leq \frac{\varepsilon}{2}.$$

By the assumption and the bounded convergence theorem, for large  $n$ ,

$$\frac{1}{\delta} \int_{-\delta}^{\delta} |\varphi(t) - \varphi_n(t)| dt \leq \frac{\varepsilon}{2}.$$

It follows that

$$\frac{1}{\delta} \int_{-\delta}^{\delta} (1 - \varphi_n(t)) dt \leq \varepsilon.$$

Hence

$$\mu_n(|x| \geq 2/\delta) \leq \varepsilon.$$

It follows that  $\mu_n$  is tight, the proof is completed.  $\square$

**COROLLARY 8.** *Let  $\mu_n$  be probability measures with characteristic functions  $\varphi_n$ . If*

$$\lim_n \varphi_n(t) = g(t), \quad \forall t,$$

*and  $g(t)$  is continuous at  $t = 0$ . Then there exists a probability measure  $\mu$  so that*

$$\mu_n \rightarrow \mu \text{ weakly,}$$

*and  $g$  is the characteristic function of  $\mu$ .*

**PROOF.** A close inspection of the proof of Theorem 117 indicates that the continuity of  $g$  at  $t = 0$  ensures the tightness of  $\{\mu_n\}$ , and at this stage  $g$  does not have to be a characteristic function. Note

$\varphi_n(0) = 1$ , hence necessarily  $g(0) = 1$ . Now knowing that  $\{\mu_n\}$  is tight, the same reasoning as before shows  $\mu_n$  converges weakly to some probability measure  $\mu$  and  $g$  is actually the characteristic function of  $\mu$ .  $\square$

### 13. Central limit theorem – the proof

LEMMA 118. *Let  $z_1, \dots, z_n, w_1, \dots, w_n$  be complex numbers of modulus  $\leq 1$ . Then*

$$\left| \prod_{k=1}^n z_k - \prod_{k=1}^n w_k \right| \leq \sum_{k=1}^n |z_k - w_k|.$$

PROOF. Note

$$\prod_{k=1}^n z_k - \prod_{k=1}^n w_k = (z_1 - w_1) \prod_{k=2}^n z_k + w_1 \left( \prod_{k=2}^n z_k - \prod_{k=2}^n w_k \right).$$

Hence

$$\left| \prod_{k=1}^n z_k - \prod_{k=1}^n w_k \right| \leq |z_1 - w_1| + \left| \prod_{k=2}^n z_k - \prod_{k=2}^n w_k \right|.$$

The conclusion follows by induction.  $\square$

The inequality may be viewed as a generalization of the geometry: if  $z_1, z_2, w_1, w_2$  be positive real numbers of absolute value  $\leq 1$ , then the difference between the areas of the rectangles  $[0, z_1] \times [0, z_2]$  and  $[0, w_1] \times [0, w_2]$  is bounded from above by the sum

$$|z_1 - w_1| \cdot 1 + 1 \cdot |z_2 - w_2|.$$

LEMMA 119. *If  $c_n \rightarrow c \in \mathbb{C}$ , then*

$$\left(1 + \frac{c_n}{n}\right)^n \rightarrow e^c.$$

PROOF. Recalling the expansion

$$e^z = \sum_{k=0}^{\infty} \frac{z^k}{k!}, \quad z \in \mathbb{C},$$

we have

$$\begin{aligned} \left| \left(1 + \frac{c_n}{n}\right)^n - e^c \right| &= \left| \prod_{k=1}^n \left(1 + \frac{c_n}{n}\right) - \prod_{k=1}^n e^{c/n} \right| \leq \sum_{k=1}^n \left| 1 + \frac{c_n}{n} - e^{c/n} \right| \\ &\leq \sum_{k=1}^n \left| 1 + \frac{c_n}{n} - \left(1 + \frac{c}{n} + o(n^{-1})\right) \right| \leq |c_n - c| + o(1). \end{aligned}$$

□

**Thm 120 (I.i.d. sequence).** *Let  $X_1, \dots, X_n$  be i.i.d. with  $EX_1 = \mu$ ,  $\text{Var}(X_1) = \sigma^2 > 0$ . If  $S = X_1 + \dots + X_n$ , then*

$$\frac{S_n - n\mu}{n^{1/2}\sigma} \rightarrow \mathcal{N}(0, 1) \text{ weakly.}$$

PROOF. By considering  $X'_k = X_k - \mu$ , we can assume w.l.g. that  $\mu = 0$ , hence  $EX^2 = \sigma^2$ . From Theorem 113, we see that the characteristic function  $\varphi$  of  $X_k$  has the expansion,

$$\varphi(t) = 1 - \frac{\sigma^2 t^2}{2} + o(t^2).$$

Using Theorem 107 and 106,  $S_n/(n^{1/2}\sigma)$  has characteristic function

$$\left(\varphi\left(\frac{t}{n^{1/2}\sigma}\right)\right)^n = \left(1 - \frac{\sigma^2}{2} \frac{t^2}{n\sigma^2} + o\left(\frac{t^2}{n\sigma^2}\right)\right)^n \rightarrow e^{-t^2/2}.$$

Note  $t$  is considered fixed relative to  $n$ . The limit is the characteristic function of  $\mathcal{N}(0, 1)$ . Whence Theorem 117 gives the desired conclusion.  $\square$

**Example 50 (Normal approximation to the binomial).** *Let  $X_1, \dots, X_n$  be independent with*

$$P(X_1 = 1) = p, \quad P(X_1 = 0) = 1 - p, \quad 0 < p < 1.$$

$S_n = X_1 + \dots + X_n$ . Then the central limit theorem tells us that,  $\forall a, b$ ,

$$\lim_{n \rightarrow \infty} P\left(a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b\right) \rightarrow \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx.$$

Imagine repeatedly tossing a fair coin and that  $X_1, \dots, X_n$  are the outcomes. We want to find the probability that in  $n = 10000$  tosses,



the number of heads is between 4900 and 5100,

$$P(4900 \leq S_n \leq 5100) = \sum_{k=4900}^{5100} C_{10000}^k \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{n-k}.$$

A direct calculation of the sum would be difficult. But using the normal approximation to the binomial, we can have a descent approximation. Note

$$ES_n = \frac{n}{2} = 5000, \quad \sqrt{\text{Var}(S_n)} = \frac{\sqrt{n}}{2} = 50.$$

Then  $S_n$  can be approximated by  $\mathcal{N}(5000, 50^2)$ ,

$$\begin{aligned} P(4900 \leq S_n \leq 5100) &\approx \frac{1}{\sqrt{2\pi} \cdot 50} \int_{4900}^{5100} e^{-\frac{1}{2} \left(\frac{x-5000}{50}\right)^2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-2}^2 e^{-\frac{x^2}{2}} dx. \end{aligned}$$

The table of standard normal distribution tells us that the integral is 0.9546.

In the example, the integer-valued random variable  $S_n$  is approximated by a continuous distribution. For each integer  $x$ , the probability  $P(S_n = x)$  is represented by the area of a rectangle with the base  $[x - 1/2, x + 1/2]$  and height  $P(X = x)$ . From this we see that for any integers  $a, b$ , the probability  $P(a \leq S_n \leq b)$  is the sum of areas of the rectangles centered at  $a, a + 1, \dots, b$ . This motivates the common practice that adjusts the integral in the example as follows,

$$P(4900 \leq S_n \leq 5100) \approx \frac{1}{\sqrt{2\pi} \cdot 50} \int_{4900-1/2}^{5100+1/2} e^{-\frac{1}{2}\left(\frac{x-5000}{50}\right)^2} dx.$$

The adjustment is called the **correction for continuity**.

**Example 51 (Normal approximation to the Poisson).** *Let  $\lambda_1, \lambda_2 > 0$ . If  $X \sim \text{Poisson}(\lambda_1)$ ,  $Y \sim \text{Poisson}(\lambda_2)$  are independent, then  $X + Y \sim \text{Poisson}(\lambda_1 + \lambda_2)$ . Thus, given  $X_1, \dots, X_n$  i.i.d.  $\sim \text{Poisson}(1)$ , the sum  $N_n = X_1 + \dots + X_n$  has*

$$N_n \sim \text{Poisson}(n).$$

Since  $EN_n = \text{Var}(N_n) = n$ , the central limit theorem yields

$$\frac{N_n - n}{\sqrt{n}} \rightarrow \mathcal{N}(0, 1) \text{ weakly.}$$

This provides an approximation of  $\text{Poisson}(\lambda)$  when  $\lambda$  is a large integer. For Poisson random variable  $N_\lambda$  with non-integer parameter  $\lambda > 0$ , let  $m = [\lambda]$ ,  $\delta = \lambda - m$ . Consider the independent random variables constructed on a common probability space (so that it is legal to talk about events that make comparison between them),

$$N_m \sim \text{Poisson}(m), \quad N_\delta \sim \text{Poisson}(\delta), \quad N_{1-\delta} \sim \text{Poisson}(1 - \delta),$$

then

$$N_\lambda \stackrel{d}{=} N_m + N_\delta, \quad N_{m+1} \stackrel{d}{=} N_m + N_\delta + N_{1-\delta}.$$

Then (noting  $N_\delta, N_{1-\delta} \geq 0$ )

$$N_m \leq N_\lambda \leq N_{m+1}.$$

Hence for  $\varepsilon > 0$  we have, as soon as  $n$  is large,

$$\frac{N_m - m}{\sqrt{m}} - \varepsilon \leq \frac{N_\lambda - \lambda}{\sqrt{\lambda}} \leq \frac{N_{m+1} - (m+1)}{\sqrt{m+1}} + \varepsilon.$$

Therefore for  $b \in \mathbb{R}$ ,

$$P\left(\frac{N_{m+1} - (m+1)}{\sqrt{m+1}} \leq b - \varepsilon\right) \leq P\left(\frac{N_\lambda - \lambda}{\sqrt{\lambda}} \leq b\right) \leq P\left(\frac{N_m - m}{\sqrt{m}} \leq b + \varepsilon\right).$$

Letting  $m \rightarrow \infty$  followed by  $\varepsilon \rightarrow 0$  shows that

$$\frac{N_\lambda - \lambda}{\sqrt{\lambda}} \rightarrow \mathcal{N}(0, 1) \text{ weakly as } \lambda \rightarrow \infty.$$

We have prove the strong law in Theorem 89 under the pairwise independence assumption. But unfortunately, the independence cannot be weakened for central limit theorem.

**Example 52 (Pairwise independence is not enough).** Let  $\xi_1, \dots, \xi_n, \dots$  be independent with

$$P(\xi_1 = 1) = P(\xi_1 = -1) = 1/2.$$

*Define successively*

$$X_1 = \xi_1; \quad (X_{2^{n-1}+1}, \dots, X_{2^{n-1}+2^{n-1}}) = \xi_{n+1} \cdot (X_1, \dots, X_{2^{n-1}}), \quad n \geq 1.$$

*Let  $S_n = X_1 + \dots + X_n$ . Then*

$$S_1 = \xi_1; \quad S_{2^n} = S_{2^{n-1}} + \xi_{n+1} \cdot S_{2^{n-1}} = S_{2^{n-1}}(1 + \xi_{n+1}), \quad n \geq 1.$$

*Thus we have by induction*

$$S_{2^n} = \xi_1(1 + \xi_2) \cdots (1 + \xi_{n+1}).$$

*Clearly*

$$P(S_{2^n} = i) = \begin{cases} 1/2^{n+1}, & i = -2^n, \\ 1/2^{n+1}, & i = 2^n, \\ 1 - 1/2^n, & i = 0. \end{cases}$$

*Hence*

$$ES_{2^n} = 0, \quad \text{Var}(S_{2^n}) = 2^n.$$

*Note each  $X_k$  takes values in  $\{-1, 1\}$  and  $EX_k = 0$ . Moreover each  $X_k$  is a product of a different finite set of  $\xi_i$ 's, and in  $X_j X_k = 0$ ,  $j \neq k$ ,*

at least one  $\xi_i$  appears exactly once. Hence it is verifiable that

$$EX_j X_k = 0, \quad \forall j \neq k.$$

Thus  $X_1, \dots, X_n, \dots$  are pairwise independent and clearly the central limit theorem does not hold for the subsequence  $n \mapsto S_{2^n}$ .

**Thm 121 (The Lindeberg-Feller theorem).** For each  $n$ , let

$X_{n,k}$ ,  $1 \leq k \leq n$  be independent,  $S_n = \sum_{k=1}^n X_{n,k}$  and

$$EX_{n,k} = 0; \quad \sigma_{n,k}^2 = EX_{n,k}^2; \quad s_n^2 = \sum_{k=1}^n \sigma_{n,k}^2.$$

Suppose that  $s_n > 0$  for large  $n$  and the **Lindeberg condition** is satisfied,

$$\forall \varepsilon > 0, \quad \lim_n \frac{1}{s_n^2} \sum_{k=1}^n E \left( X_{n,k}^2 1_{\{|X_{n,k}| \geq \varepsilon s_n\}} \right) = 0.$$

Then

$$\frac{S_n}{s_n} \rightarrow \mathcal{N}(0, 1) \text{ weakly.}$$

PROOF. By replacing  $X_{n,k}$  with  $X_{n,k}/s_n$ , we may assume w.l.g. that

$$s_n^2 = \sum_{k=1}^n \sigma_{n,k}^2 = 1.$$

Write  $\varphi_{n,k}$  for the characteristic function of  $X_{n,k}$ . By the continuity theorem (Theorem 117), it is sufficient to show that

$$\prod_{k=1}^n \varphi_{n,k} \rightarrow e^{-\frac{1}{2}t^2}, \quad \forall t \in \mathbb{R}.$$

But this is equivalent to

$$\left| \prod_{k=1}^n \varphi_{n,k} - \prod_{k=1}^n e^{-\frac{1}{2}\sigma_{n,k}^2 t^2} \right| \rightarrow 0, \quad \forall t \in \mathbb{R}.$$

(Note  $s_n^2 = 1$ ). In view of Lemma 118, we have

$$(13.1) \quad \left| \prod_{k=1}^n \varphi_{n,k} - \prod_{k=1}^n e^{-\frac{1}{2}\sigma_{n,k}^2 t^2} \right| \leq \sum_{k=1}^n \left| \varphi_{n,k} - e^{-\frac{1}{2}\sigma_{n,k}^2 t^2} \right|.$$

But the RHS is less than the sum of

$$\sum_{k=1}^n \left| \varphi_{n,k} - \left( 1 - \frac{1}{2}t^2\sigma_{n,k}^2 \right) \right| \text{ and } \sum_{k=1}^n \left| e^{-\frac{1}{2}\sigma_{n,k}^2 t^2} - \left( 1 - \frac{1}{2}t^2\sigma_{n,k}^2 \right) \right|.$$

Therefore it remains to prove that the sums converge to zero as  $n \rightarrow \infty$ . Since  $EX_{n,k}^2 < \infty$  by the assumption, we get from Lemma 113 that

$$\left| \varphi_{n,k} - \left( 1 - \frac{1}{2}t^2\sigma_{n,k}^2 \right) \right| \leq E \min\{|tX_{n,k}|^2, |tX_{n,k}|^3\}.$$



For any  $\varepsilon > 0$ ,

$$\begin{aligned}
E \min\{|tX_{n,k}|^2, |tX_{n,k}|^3\} &= \int_{|X_{n,k}| < \varepsilon} \min\{\cdot, \cdot\} dP + \int_{|X_{n,k}| \geq \varepsilon} \min\{\cdot, \cdot\} dP \\
&\leq \int_{|X_{n,k}| < \varepsilon} |tX_{n,k}|^3 dP + \int_{|X_{n,k}| \geq \varepsilon} |tX_{n,k}|^2 dP \\
&\leq \varepsilon |t|^3 \sigma_{n,k}^2 + t^2 \int_{|X_{n,k}| \geq \varepsilon} X_{n,k}^2 dP.
\end{aligned}$$

Hence

$$\sum_{k=1}^n \left| \varphi_{n,k} - \left( 1 - \frac{1}{2} t^2 \sigma_{n,k}^2 \right) \right| \leq \varepsilon |t|^3 + t^2 \sum_{k=1}^n \int_{|X_{n,k}| \geq \varepsilon} X_{n,k}^2 dP.$$

By the arbitrariness of  $\varepsilon$  and the Lindeberg condition, this gives

$$(13.2) \quad \lim_n \sum_{k=1}^n \left| \varphi_{n,k} - \left( 1 - \frac{1}{2} t^2 \sigma_{n,k}^2 \right) \right| = 0.$$

Since

$$\sigma_{n,k}^2 = EX_{n,k}^2 \leq \varepsilon^2 + \int_{|X_{n,k}| \geq \varepsilon} X_{n,k}^2 dP,$$

so the Lindeberg condition implies that

$$\max_{1 \leq k \leq n} \sigma_{n,k}^2 \rightarrow 0.$$

Note

$$|e^z - 1 - z| \leq |z|^2 e^{|z|}, \quad \forall z \in \mathbb{C}.$$

Hence

$$\sum_{k=1}^n \left| e^{-\frac{1}{2}\sigma_{n,k}^2 t^2} - \left( 1 - \frac{1}{2}t^2 \sigma_{n,k}^2 \right) \right| \leq t^4 \sum_{k=1}^n \sigma_{n,k}^4 e^{t^2} \rightarrow 0.$$

Therefore, this together with (13.2) shows that the RHS of (13.1) converges to zero as  $n \rightarrow \infty$ . The proof is completed.  $\square$

**Example 53 (Random permutation and record values).**  
*Continuing Example 39 and Example 41, let  $Y_1, \dots, Y_n$  be independent*

with

$$P(Y_k = 1) = \frac{1}{k}, \quad P(Y_k = 0) = 1 - \frac{1}{k}.$$

Let  $S_n = Y_1 + \cdots + Y_n$ . Since

$$EY_k = \frac{1}{k}, \quad \text{Var}(Y_k) = \frac{1}{k} - \frac{1}{k^2},$$

we have

$$ES_n \sim \log n, \quad \text{Var}(S_n) \sim \log n.$$

To invoke the central limit theorem on  $S_n$ , consider the triangular array

$$X_{n,k} = \frac{Y_k - 1/k}{\sqrt{\log n}}, \quad 1 \leq k \leq n.$$

Clearly

$$EX_{n,k} = 0, \quad \sum_{k=1}^n X_{n,k}^2 \rightarrow 1 \text{ as } n \rightarrow \infty.$$

Since

$$|X_{n,k}| \leq \frac{1}{\sqrt{\log n}}, \quad 1 \leq k \leq n,$$

the Lindeberg condition is satisfied:  $\forall \varepsilon > 0$ ,

$$\sum_{k=1}^n E\left(X_{n,k}^2 1_{\{|X_{n,k}| \geq \varepsilon\}}\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

So Theorem 121 shows that

$$\frac{S_n - \sum_{k=1}^n 1/k}{\sqrt{\log n}} \rightarrow \mathcal{N}(0, 1) \text{ weakly.}$$

Since  $\sum_{k=1}^n 1/k \sim \log n$ , we have

$$\frac{S_n - \log n}{\sqrt{\log n}} \rightarrow \mathcal{N}(0, 1) \text{ weakly.}$$