

# Appendix B

## Mathematical Background on Probabilities

Po-Ning Chen, Professor

Department of Communications Engineering

National Chiao Tung University

Hsin Chu, Taiwan 300, R.O.C.

# Probability space and random variable

I: b-1

- A *probability space* is a triple  $(\Omega, \mathcal{F}, P)$ , where
  - $\Omega$  is the set of all possible outcomes (often named *sample space*), and
  - $\mathcal{F}$  is the  $\sigma$ -field of  $\Omega$  (often named *event space*), and
  - $P$  is a probability measure on the  $\sigma$ -field, satisfying:
    1.  $0 \leq P(A) \leq 1$  for  $A \in \mathcal{F}$ ;
    2.  $P(\emptyset) = 0$  and  $P(\Omega) = 1$ .
    3. **(countable additivity)** if  $A_1, A_2, \dots$  is a disjoint sequence of sets in  $\mathcal{F}$ , then

$$P\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} P(A_k).$$

- A random variable  $X$  on a probability space  $(\Omega, \mathcal{F}, P)$  is a *real-valued* function on  $\Omega$  (i.e.,  $X : \Omega \rightarrow \mathfrak{R}$ ), satisfying that the set  $\{\omega : X(\omega) = x\} \in \mathcal{F}$  for each real  $x$ .

## Probability space and random variable

I: b-2

- The observation probability space  $(\Omega_X, \mathcal{F}_X, P_X)$  of a random variable  $X$  defined over the probability space  $(\Omega, \mathcal{F}, P)$  satisfies:

$$\begin{cases} \Omega_X = X(\Omega) \\ \mathcal{F}_X = \{X(A) \subset \mathfrak{R} : A \in \mathcal{F}\} \\ P_X(G) = P(\{\omega \in \Omega : X(\omega) \in G\}) \text{ for any } G \in \mathcal{F}_X, \end{cases}$$

where  $X(A) = \{x \in \mathfrak{R} : X(\omega) = x \text{ for some } \omega \in A\}$ .

## Why define random variables based on $(\Omega, \mathcal{F}, P)$ ?

I: b-3

Answer:  $(\Omega, \mathcal{F}, P)$  is what truly occurs internally, but possibly **non-observable**.

- In order to infer which of the *non-observable*  $\omega$  occurs, an experiment that results in observable  $x$  that is a function of  $\omega$  is performed.
- So  $x$  is a function of  $\omega$ , which takes real values.
- Such an experiment results in the random variable  $X$  whose probability is so defined over the probability space  $(\Omega, \mathcal{F}, P)$ .

## Why define random variables based on $(\Omega, \mathcal{F}, P)$ ?

---

I: b-4

- In applications, we are more interested in the observation probability space  $(\Omega_X, \mathcal{F}_X, P_X)$  than the inherited probability space  $(\Omega, \mathcal{F}, P)$  on which the random variable is defined.
- It can be proved that given a real-valued non-negative function  $F(\cdot)$ , satisfying that  $\lim_{x \downarrow -\infty} F(x) = 0$  and  $\lim_{x \uparrow \infty} F(x) = 1$ , there exist a random variable and an inherited probability space such that the cumulate distribution function (cdf) of the random variable defined over the probability space is equal to  $F(\cdot)$ .
- The above result releases us with the burden of referring to a probability space before our defining a random variable. In other words, we can indeed define a random variable  $X$  directly by its cdf, i.e.,  $\Pr[X \leq x]$ , without bothering to refer to its inherited probability space. Nevertheless, it is better to keep in mind (and learn) that a formal mathematical notion of random variables is defined over some probability space.
- In what follows, you will notice that most of the properties of random variables (random processes) are defined based on their *observation probability space*.

## Random Processes

I: b-5

- A random process  $\mathbf{X} = \{X_t, t \in I\}$  is a collection of random variables that arise from *common* probability space  $(\Omega, \mathcal{F}, P)$ .
  - Under this definition, any finite dimensional distribution of  $\{X_t, t \in I\}$  is well-defined. For example,

$$\begin{aligned} & \Pr[X_1 \leq x_1, X_5 \leq x_5, X_9 \leq x_9] \\ &= P(\{\omega \in \Omega : X_1(\omega) \leq x_1, X_5(\omega) \leq x_5, X_9(\omega) \leq x_9\}) \end{aligned}$$

- Again, we only consider its distribution  $P_{\mathbf{X}}$  and *observation event space*  $\mathcal{F}_{\mathbf{X}}$  in the sequel.

## Random Processes in Communications

I: b-6

$$\mathbf{X} = \dots, X_{-3}, X_{-2}, X_{-1}, X_0, X_1, X_2, X_3, \dots$$

- Memoryless = independent and identically distributed (i.i.d.)
- First-order stationary =  $(\forall i, j \in I) P_{X_i} = P_{X_j}$
- Second-order stationary =  $(\forall i, j \in I \text{ and } \forall k) P_{X_i, X_j} = P_{X_{i+k}, X_{j+k}}$
- (Strictly) stationary =  $k$ -th-order stationary for all  $k$
- Weakly stationary

$$(\forall i, j \in I \text{ and } \forall k) \quad E[X_i] = E[X_j] \text{ and } E[X_i X_j] = E[X_{i+k} X_{j+k}]$$

## Random Processes in Communications

I: b-7

$$\mathbf{X} = \dots, X_{-3}, X_{-2}, X_{-1}, X_0, X_1, X_2, X_3, \dots$$

- Ergodic

- Time-shift invariant subset/event in  $\mathcal{F}_X$ :
  - \* Recall that an event is a subset of the sample space  $\mathbf{X}(\Omega)$ , and is an element of the observation event space  $\mathcal{F}_X$ .
  - \* A time-shift invariance event  $A \in \mathcal{F}_X$  is one that the time shift counterpart of any element in  $A$  still belong to  $A$ . E.g.,

$$A = \{(\dots, x_{-1} = 0, x_0 = 1, x_1 = 0, x_2 = 1, \dots), \\ (\dots, x_{-1} = 1, x_0 = 0, x_1 = 1, x_2 = 0, \dots)\}.$$

- For any two time-shift invariant events, either they are disjoint, or one is contained in the other, or their intersection is also time-shift invariant.
- Ergodic = every time-shift invariant event has either probability one or zero!



# Random Processes in Communications

I: b-8

- Ergodic = every time-shift invariant event (in  $\mathcal{F}_{\mathbf{X}}$ ) has either probability one or zero.

– *Interpretations:* At most one of the following time-shift invariant events has probability one.

$$\left\{ x_{-\infty}^{\infty} \in \{0, 1\}_{-\infty}^{\infty} : \lim_{n \rightarrow \infty} \frac{x_{-n} + \dots + x_0 + \dots + x_n}{2n + 1} = 0.0 \right\}$$

$$\left\{ x_{-\infty}^{\infty} \in \{0, 1\}_{-\infty}^{\infty} : \lim_{n \rightarrow \infty} \frac{x_{-n} + \dots + x_0 + \dots + x_n}{2n + 1} = 0.1 \right\}$$

⋮

$$\left\{ x_{-\infty}^{\infty} \in \{0, 1\}_{-\infty}^{\infty} : \lim_{n \rightarrow \infty} \frac{x_{-n} + \dots + x_0 + \dots + x_n}{2n + 1} = 0.9 \right\}$$

$$\left\{ x_{-\infty}^{\infty} \in \{0, 1\}_{-\infty}^{\infty} : \lim_{n \rightarrow \infty} \frac{x_{-n} + \dots + x_0 + \dots + x_n}{2n + 1} = 1.0 \right\}$$

★ Under boundedness, an outcome has the same limit (if it exists) as its shift counterpart.

$$\left| \frac{x_{-n} + \dots + x_n}{2n + 1} - \frac{x_{-n+1} + \dots + x_{n+1}}{2n + 1} \right| = \left| \frac{x_{-n} - x_{n+1}}{2n + 1} \right| \leq \frac{1}{2n + 1}.$$

# Random Processes in Communications

I: b-9

- Stationarity  $\Rightarrow$  Convergence of time average to a random variable
  - (One-sided) *Time average* is what we can truly observe and calculate!
  - E.g.,  $\{X_i\}_{i=1}^n$  i.i.d. with  $\Pr\{X_i = 1\} = 1 - \Pr\{X_i = 0\} = U$ , where  $\Pr\{U = 0.2\} = \Pr\{U = 0.8\} = 0.5$ .

$$\Rightarrow \frac{X_1 + X_2 + \dots + X_n}{n} \rightarrow U.$$

- Ergodicity and boundedness  $\Rightarrow$  Time average converges to a constant

– It does not guarantee that the *constant* is exactly the *ensemble mean*.

- E.g.  $\Pr\{(\dots, x_{-1} = 0, x_0 = 1, x_1 = 0, x_2 = 1, \dots)\} = 0.2$   
and  $\Pr\{(\dots, x_{-1} = 1, x_0 = 0, x_1 = 1, x_2 = 0, \dots)\} = 0.8$ .

Then

$$\frac{X_{-n} + \dots + X_0 + \dots + X_n}{2n + 1} \rightarrow \frac{1}{2};$$

but  $E[X_i] =$  either 0.8 or 0.2.

- **Note:** Two alternative names for *time average* are *sample average* and *Cesàro mean*. Throughout the slides, we will use the *time average* since it gives more intuition to its formula.

# Random Processes in Communications

I: b-10

- Stationary ergodicity  $\Rightarrow$  Convergence of time average to ensemble average.

**In mathematics:**

$$\text{Stationarity} \Rightarrow \Pr \left\{ \lim_{n \rightarrow \infty} \frac{X_{-n} + \dots + X_0 + \dots + X_n}{2n + 1} = Y \right\} = 1$$

for some random variable  $Y$

$$\text{Ergodicity and boundedness} \Rightarrow \Pr \left\{ \lim_{n \rightarrow \infty} \frac{X_{-n} + \dots + X_0 + \dots + X_n}{2n + 1} = a \right\} = 1$$

for some scalar  $a$

$$\text{Stationary Ergodicity} \Rightarrow \Pr \left\{ \lim_{n \rightarrow \infty} \frac{X_{-n} + \dots + X_0 + \dots + X_n}{2n + 1} = E[X_0] \right\} = 1$$

# Random Processes in Communications

I: b-11

## Physical meaning of stationary ergodic assumption

- Stationary ergodic random source
  - Empirical distribution (relative frequency) can be used to approximate the true distribution.
  - Note that empirical distribution is what we can truly observe and calculate! E.g., observe a dice rolling experiment and obtain

$$X_1^{30} = 154326543334225632425644234443$$

Given that the experimental source is stationary ergodic, we can approximate the true distribution by (cf. pointwise ergodic theorem or Theorem B.1):

$$\begin{array}{lll} \Pr\{X_i = 1\} & \approx \frac{1}{30} & \Pr\{X_i = 2\} \approx \frac{6}{30} & \Pr\{X_i = 3\} \approx \frac{7}{30} \\ \Pr\{X_i = 4\} & \approx \frac{9}{30} & \Pr\{X_i = 5\} \approx \frac{4}{30} & \Pr\{X_i = 6\} \approx \frac{3}{30} \end{array}$$

- Non-stationary or non-ergodic source
  - Empirical distribution (relative frequency) cannot necessarily be used to approximate the true distribution.

## Random Processes in Communications

I: b-12

### **Physical meaning of stationary ergodic assumption**

- In theories of communications, people assume that *the source is stationary* or *the source is stationary ergodic*. But you seldom find the assumption of *the source being ergodic but non-stationary*. Why?
  - Because an ergodic but non-stationary source not only does not facilitate the analytical study of communications problems, but seems no applications in practice.
  - From this, we learn that assumptions are made either to ease our analytical study of communications problem or to fit a specific need of applications. Without the two footings, an assumption becomes of minor interest.
- In other words, *ergodicity* assumption usually comes after *stationarity* assumption. A specific example is the **pointwise ergodic theorem**, where the random processes considered is presumed to be **stationary**.

## Random Processes in Communications

I: b-13

**Theorem B.1 (pointwise ergodic theorem)** Give a discrete-time stationary random process  $\{X_n\}_{-\infty < n < \infty}$ . For arbitrary real-valued function  $f(\cdot)$  on  $\mathfrak{R}$  with finite mean ( $|E[f(X_n)]| < \infty$ ), there exists a random variable  $\hat{Y}$  such that

$$\Pr \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{n=0}^{\infty} f(X_n) = \hat{Y} \right] = 1.$$

If, in addition to stationarity, the process is also ergodic, then

$$\Pr \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{n=0}^{\infty} f(X_n) = E[\hat{Y}] \right] = 1.$$

**Example B.2** Consider the process  $\{X_i\}_{i=-\infty}^{\infty}$  consisting of a family of i.i.d. binary random variables (obviously, it is stationary and ergodic). Define the function  $f(\cdot)$  by  $f(0) = 0$  and  $f(1) = 1$ , Hence,

$$E[f(X)] = P_X(0)f(0) + P_X(1)f(1) = P_X(1)$$

is finite. By the pointwise ergodic theorem, we have

$$\lim_{n \rightarrow \infty} \frac{f(X_1) + f(X_2) + \dots + f(X_n)}{n} = \lim_{n \rightarrow \infty} \frac{X_1 + X_2 + \dots + X_n}{n} = P_X(1).$$

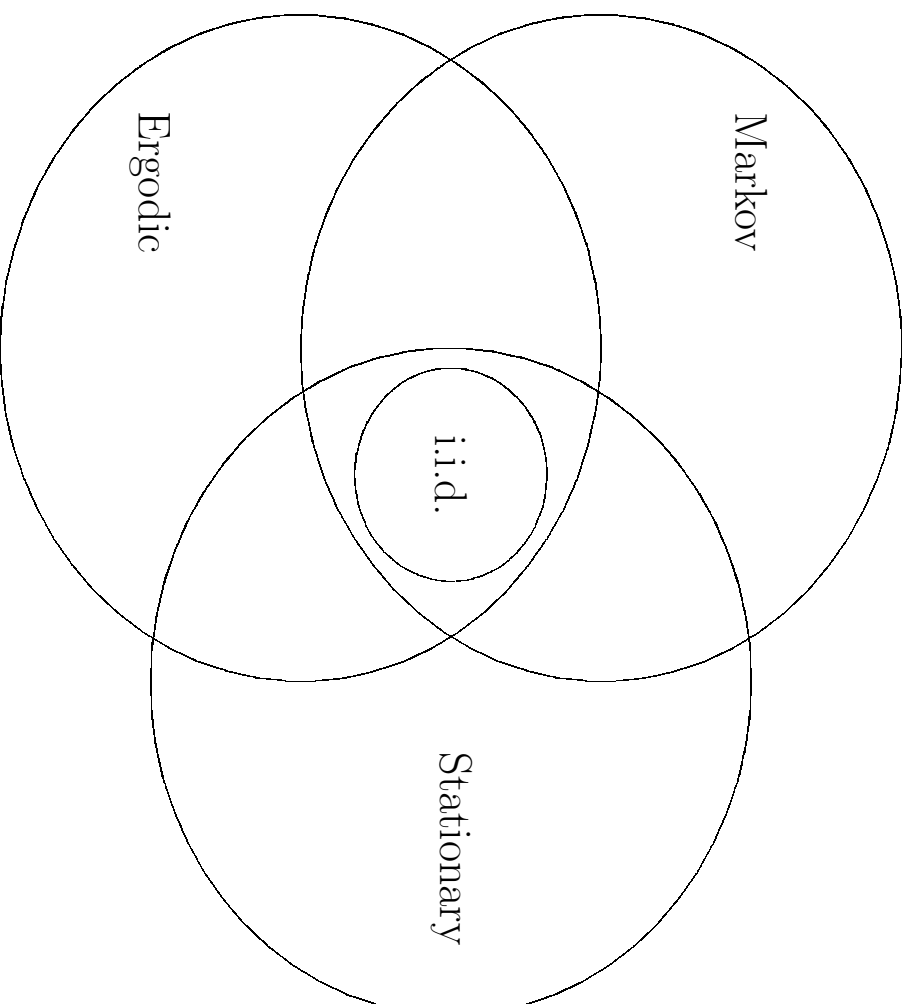
# Random Processes in Communications

I: b-14

- First-order Markov: Denote by  $X \rightarrow Y \rightarrow Z$  or  $Z \rightarrow Y \rightarrow X$  or  $X \leftrightarrow Y \leftrightarrow Z$ .
  - In words,  $X$  and  $Z$  are conditionally independent given  $Y$
- $k$ th-order Markov
$$\Pr\{X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_1 = x_1\}$$
$$= \Pr\{X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_{n-k} = x_{n-k}\}.$$
- ( $k$ -th order) Irreducible Markov
  - No reduction on the number of possible state outcome is rendered.
- Homogeneous or time-invariant Markov
$$(\forall x^k, y^k \text{ and } j) \quad \Pr\{X_j^{k+j-1} = x^k | X_1^k = y^k\} > 0.$$
- The transition probability is time-invariant.
- Stationary Markov = Homogeneous with stationary initial distribution.

# General Relation of Random Processes

I: b-15





# Convergences of Sequences of Random Variables

I: b-16

- Relation of five modes of convergence

$$\begin{array}{ccccc} X_n & \xrightarrow{p.w.} & X & & \\ \Downarrow & & & & \\ X_n & \xrightarrow{a.s.} & X & \xrightarrow[\text{Thm. B.8}]{\text{Thm. B.7}} & X_n & \xrightarrow{L^r} & X \quad (r \geq 1) \\ \Downarrow & & & & \Downarrow & & \\ X_n & \xrightarrow{p} & X & & & & \\ \Downarrow & & & & & & \\ X_n & \xrightarrow{d} & X & & & & \end{array}$$

# Convergences of Sequences of Random Variables

I: b-17

- Pointwise convergence and almost surely convergence

**E.g.** Give a probability space

$$(\Omega = \{0, 1, 2, 3\}, 2^\Omega, P(0) = P(1) = P(2) = 1/3).$$

- A random variable  $X_n$  is a mapping from a probability space to  $\mathfrak{R}$ . Let the mapping be

$$X_n(\omega) = \frac{\omega}{n} \Rightarrow \Pr\{X_n = 0\} = \Pr\left\{X_n = \frac{1}{n}\right\} = \Pr\left\{X_n = \frac{2}{n}\right\} = \frac{1}{3}.$$

- (*Pointwise convergence*) Observe that

$$(\forall \omega \in \Omega) X_n(\omega) \rightarrow X(\omega),$$

where  $X(\omega) = 0$  for every  $\omega \in \Omega$ . So

$$X_n \xrightarrow{p.w.} X.$$

- (*Almost surely convergence*) Let  $\tilde{X}(\omega) = 0$  for  $\omega = 0, 1, 2$  and  $\tilde{X}(\omega) = 1$  for  $\omega = 3$ . Then both of the following statements are true:

$$X_n \xrightarrow{a.s.} X \quad \text{and} \quad X_n \xrightarrow{a.s.} \tilde{X},$$

(since

$$\Pr\left\{\lim_{n \rightarrow \infty} X_n = \tilde{X}\right\} = \sum_{\omega=0}^3 P(\omega) \cdot \mathbf{1}\left\{\lim_{n \rightarrow \infty} X_n(\omega) = \tilde{X}(\omega)\right\} = 1.)$$

## Convergences of Sequences of Random Variables

I: b-18

- Almost surely convergence (with probability 1) and convergence in probability

$$X_n \xrightarrow{a.s.} X \equiv \Pr \left\{ \lim_{n \rightarrow \infty} X_n = X \right\} = 1$$

$$X_n \xrightarrow{p} X \equiv (\forall \gamma > 0) \lim_{n \rightarrow \infty} \Pr \{ |X_n - X| < \gamma \} = 1$$

- Convergence in  $r$ th mean

$$X_n \xrightarrow{L_r} X \equiv \lim_{n \rightarrow \infty} E [|X_n - X|^r] = 0$$

- Convergence in distribution

$$X_n \xrightarrow{d} X \equiv \lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x) \text{ for every continuous point of } F_X(x)$$

# Convergences of Sequences of Random Variables

I: b-19

The next observation facilitates the finding of limiting random variable.

## **Observation B.5 (uniqueness of convergence)**

1. If  $X_n \xrightarrow{p.w.} X$  and  $X_n \xrightarrow{p.w.} Y$ ,  
then  $X = Y$  pointwisely. I.e.,

$$(\forall \omega \in \Omega) \quad X(\omega) = Y(\omega).$$

2. If  $X_n \xrightarrow{a.s.} X$  and  $X_n \xrightarrow{a.s.} Y$   
(or  $X_n \xrightarrow{p} X$  and  $X_n \xrightarrow{p} Y$ )  
(or  $X_n \xrightarrow{L_r} X$  and  $X_n \xrightarrow{L_r} Y$ ),  
then  $X = Y$  with probability 1. I.e.,

$$\Pr\{X = Y\} = 1.$$

3.  $X_n \xrightarrow{d} X$  and  $X_n \xrightarrow{d} Y$ ,  
then  $F_X(x) = F_Y(x)$  for every continuity point of  $F_X(x)$ .

# Convergences of Sequences of Random Variables

I: b-20

The next observation facilitates the proofs of convergence of sequences of random variables.

## **Observation B.6 (mutual convergence criteria)**

1.  $\{X_n\}_{n=1}^\infty$  converges pointwisely if, and only if,

$$(\forall \omega \in \Omega) \lim_{n \rightarrow \infty} |X_{n+1}(\omega) - X_n(\omega)| = 0.$$

2.  $\{X_n\}_{n=1}^\infty$  converges with probability 1 if, and only if,

$$\begin{aligned} P \left\{ \omega \in \Omega : \lim_{n \rightarrow \infty} |X_{n+1}(\omega) - X_n(\omega)| = 0 \right\} \\ = \Pr \left\{ \lim_{n \rightarrow \infty} |X_{n+1} - X_n| = 0 \right\} = 1. \end{aligned}$$

3.  $\{X_n\}_{n=1}^\infty$  converges in probability if, and only if, for every  $\varepsilon > 0$ ,

$$\begin{aligned} \lim_{n \rightarrow \infty} P \{ \omega \in \Omega : |X_{n+1}(\omega) - X_n(\omega)| > \varepsilon \} \\ = \lim_{n \rightarrow \infty} \Pr \{ |X_{n+1} - X_n| > \varepsilon \} = 0. \end{aligned}$$

4.  $\{X_n\}_{n=1}^\infty$  converges in  $r^{\text{th}}$  mean if, and only if,

$$\lim_{n \rightarrow \infty} E [|X_{n+1} - X_n|^r] = 0.$$

## Convergences of Sequences of Random Variables

I: b-21

The following statement is not necessarily true.

$$\boxed{\begin{array}{l} \{X_n\}_{n=1}^\infty \text{ converges in distribution if, and only if,} \\ \lim_{n \rightarrow \infty} |F_{X_{n+1}}(x) - F_{X_n}(x)| = 0 \text{ for every } x. \end{array}}$$

**E.g.**  $\Pr\{X_n = n\} = 1$  for every  $n$ . Then

$$(\forall x \in \mathcal{R}) \lim_{n \rightarrow \infty} |F_{X_{n+1}}(x) - F_{X_n}(x)| = 0.$$

But  $\{X_n\}_{n=1}^\infty$  does not converge in distribution to any random variable.

# Convergences of Sequences of Random Variables

I: b-22

**Theorem B.7 (monotone convergence theorem)**

$$\left. \begin{array}{l} (i) X_n \xrightarrow{a.s.} X \\ (ii) (\forall n) Y \leq X_n \leq X_{n+1} \\ (iii) E[|Y|] < \infty \end{array} \right\} \Rightarrow X_n \xrightarrow{L^1} X \Rightarrow E[X_n] \rightarrow E[X].$$

**Theorem B.8 (dominated convergence theorem)**

$$\left. \begin{array}{l} (i) X_n \xrightarrow{a.s.} X \\ (ii) (\forall n) |X_n| \leq Y \\ (iii) E[|Y|] < \infty \end{array} \right\} \Rightarrow X_n \xrightarrow{L^1} X \Rightarrow E[X_n] \rightarrow E[X].$$

## Law of Large Numbers

I: b-23

**Theorem B.9 (weak law of large number)** Let  $\{X_n\}_{n=1}^\infty$  be a sequence of uncorrelated random variables with common mean  $E[X_i] = \mu$ . If the variables also have common variance, or more generally,

$$\lim_{n \rightarrow \infty} \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = 0, \quad (\text{equivalently, } \frac{X_1 + \dots + X_n}{n} \xrightarrow{L^2} \mu)$$

then

$$\frac{X_1 + \dots + X_n}{n} \xrightarrow{p} \mu.$$

**proof:** By Chebyshev's inequality,

$$\Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \varepsilon \right\} \leq \frac{1}{n^2 \varepsilon^2} \sum_{i=1}^n \text{Var}[X_i].$$

□

**Note:**  $X_n \xrightarrow{L^2} X$  implies  $X_n \xrightarrow{p} X$ .



## Law of Large Numbers

I: b-24

**Theorem B.10 (Kolmogorov's strong law of large number)** Let  $\{X_n\}_{n=1}^{\infty}$  be an independent sequence of random variables with common mean  $E[X_n] = \mu$ . If either

1.  $X_n$ 's are identically distributed; or
2.  $X_n$ 's are square-integrable with

$$\sum_{i=1}^{\infty} \frac{\text{Var}[X_i]}{i^2} < \infty,$$

Then

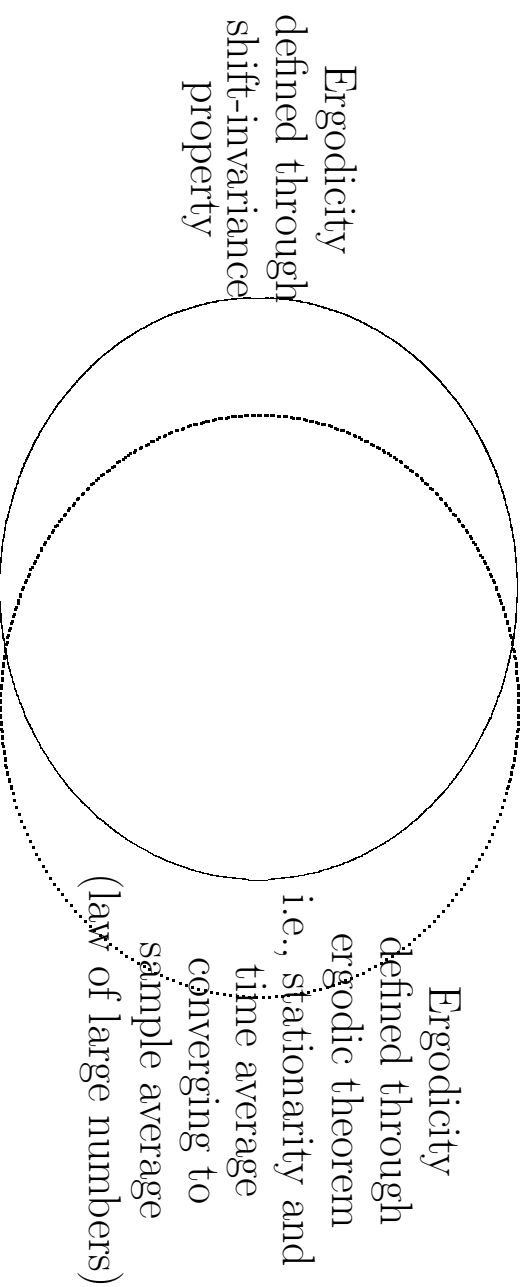
$$\frac{X_1 + \cdots + X_n}{n} \xrightarrow{a.s.} \mu.$$

**Note:** The difference of *weak* and *strong* laws of large number is that the former is convergence *in probability*, while the latter is *almost surely* convergence.

# Ergodicity and Law of Large Number

I: b-25

- Ergodicity is defined through *shift-invariant property*; ergodic theorem (with implicitly stationary assumption) is a *consequence* of this definition.
- However, some engineers define **ergodicity** as the sequences that satisfies **ergodic theorem**, which is not exactly correct in mathematics. Nevertheless, since engineers usually concern only those sequences that satisfy the *ergodic theorem* (for which case we can use time average to approximate ensemble average), such misinformed definition does not cause serious problem in its usage.



## Ergodicity and Law of Large Number

I: b-26

- The concept of ergodicity does not require stationarity. In other words, a non-stationary process can be ergodic.
- Many perfectly good models of physical processes are not ergodic, yet they have a form of law of large numbers (i.e., time average converging to ensemble average). In other words, non-ergodic processes can be perfectly good and useful models.
- There is no finite-dimensional equivalent definition of ergodicity as there is for stationarity. This fact makes it more difficult to describe and interpret ergodicity.
  - Hard to define a *shift-invariant* event for finite dimensional system except by adopting the concept of “cyclic,” which may not be apt to engineering need.
- I.i.d. processes are ergodic, i.e., ergodicity can be thought of as a (kind of) generalization of i.i.d.

## Ergodicity and Law of Large Number

I: b-27

• As mentioned earlier, stationarity and ergodicity implies the time average converges with probability 1 to the ensemble mean. Now if a process is stationary but not ergodic, then the time average still converges, but possibly not to the ensemble mean (indeed, to a random variable).

**E.g.** A stationary non-ergodic process.

- $\{A_n\}_{n=-\infty}^{\infty}$  and  $\{B_n\}_{n=-\infty}^{\infty}$  i.i.d.
- $\Pr\{A_n = 0\} = 1 - \Pr\{A_n = 1\} = \Pr\{B_n = 1\} = 1 - \Pr\{B_n = 0\} = 1/4$ .
- $V \in \{0, 1\}$ .
- $V \perp\!\!\!\perp \{A_n\}_{n=-\infty}^{\infty} \perp\!\!\!\perp \{B_n\}_{n=-\infty}^{\infty}$
- Define  $\{X_n\}_{n=-\infty}^{\infty}$  to be:

$$X_n = \begin{cases} A_n, & \text{if } V = 1; \\ B_n, & \text{if } V = 0. \end{cases}$$

– Then

$$\frac{X_1 + X_2 + \dots + X_n}{n} \xrightarrow{a.s.} U,$$

where  $\Pr\{U = 3/4\} = 1 - \Pr\{U = 1/4\} = \Pr\{V = 1\}$ .

# Ergodicity and Law of Large Number

I: b-28

- For a stationary but non-ergodic process, pointwise ergodic theorem gives the desired **limiting behavior** of the **time average**. E.g., in the previous example,

$$\frac{X_1 + X_2 + \cdots + X_n}{n} \xrightarrow{a.s.} U.$$

- Note that the sequences in the above example does not satisfy the *law of large number*.
- The previous example also shed the light to a famous and surprising result: *ergodic decomposition theorem* for which one of the main contributors is Prof. Davison at EE department of UMCP,<sup>1</sup> who is now retired.

---

<sup>1</sup>UMCP=University of Maryland, College Park

# Ergodicity and Law of Large Number

I: b-29

- Ergodic decomposition theorem
  - Under fairly general conditions, a stationary process is a mixture of stationary ergodic processes, e.g., in the previous example,

$$X_n = V \cdot A_n + (1 - V) \cdot B_n.$$

- *Implication:* One always observes a stationary ergodic outcome (or stationary-ergodic-like behavior) for any stationary(-only) process, i.e.,

if  $V = 1$ , then one observes  $\dots, A_1, A_2, A_3, \dots$

if  $V = 0$ , then one observes  $\dots, B_1, B_2, B_3, \dots$

where  $\{A_n\}_{n=-\infty}^{\infty}$  and  $\{B_n\}_{n=-\infty}^{\infty}$  are both stationary ergodic.

## Ergodicity and Law of Large Number

I: b-30

- The previous remarks hint that **ergodicity** is not required for the **strong law of large number** to be useful.
- The next question is whether or not **stationarity** is required. Again the answer is negative! In fact, the main concern of strong law of large numbers is the convergence of sample averages to its ensemble expectation. It should be reasonable to expect that random processes could exhibit transient behavior that violates the stationarity definition, yet the sample average still converges. One can then introduce the notion of *asymptotically stationary* to the strong law of large numbers. (I.e., a sequence of random variables is said to be *asymptotically stationary* if it satisfies strong law of large numbers.)
- So please do not take stationarity and ergodicity too serious. Strong- law-of-large-numbers-type behavior is really our main concern (in information theory).
- Since stationary ergodicity assumption is occasionally seen in theoretical papers, it is nonetheless better to have a basic understanding of it.

## Central Limit Theorem

I: b-31

**Theorem B.12 (central limit theorem)** If  $\{X_n\}_{n=1}^{\infty}$  is a sequence of i.i.d. random variables with common marginal mean  $\mu$  and variance  $\sigma^2$ , then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

where  $\mathcal{N}(0, \sigma^2)$  represents the Gaussian distribution with mean 0 and variance  $\sigma^2$ .



## Convexity and Concavity

I: b-32

**Definition B.13 (convexity)** A function  $f(x)$  is said to be *convex* over an interval  $(a, b)$  if for every  $x_1, x_2 \in (a, b)$  and  $0 \leq \lambda \leq 1$ ,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2).$$

Furthermore, a function  $f$  is said to be *strictly convex* if equality holds only when  $\lambda = 0$  or  $\lambda = 1$ .

**Definition B.14 (concavity)** A function  $f$  is *concave* if  $-f$  is convex.

## Jensen's inequality

I: b-33

**Theorem B.15 (Jensen's inequality)** If  $f$  is convex and  $X$  is a random variable, then

$$E[f(X)] \geq f(E[X]).$$

Moreover, if  $f$  is strictly convex, then equality in the above inequality immediately implies  $X = E[X]$  with probability 1.

**Proof:** Let  $y = ax + b$  be a support line through the point  $(E[X], f(E[X]))$ , where a support line<sup>2</sup> (for a convex function) at  $x_0$  is by definition a line passing through the point  $(x_0, f(x_0))$  and is lying entirely below the graph of  $f(\cdot)$ . Thus,

$$(\forall x \in \mathcal{X}) \quad ax + b \leq f(x).$$

By taking the expectation value of both sides, we obtain

$$a \cdot E[X] + b \leq E[f(X)],$$

but we know that  $a \cdot E[X] + b = f(E[X])$ . Consequently,

$$f(E[X]) \leq E[f(X)].$$

□

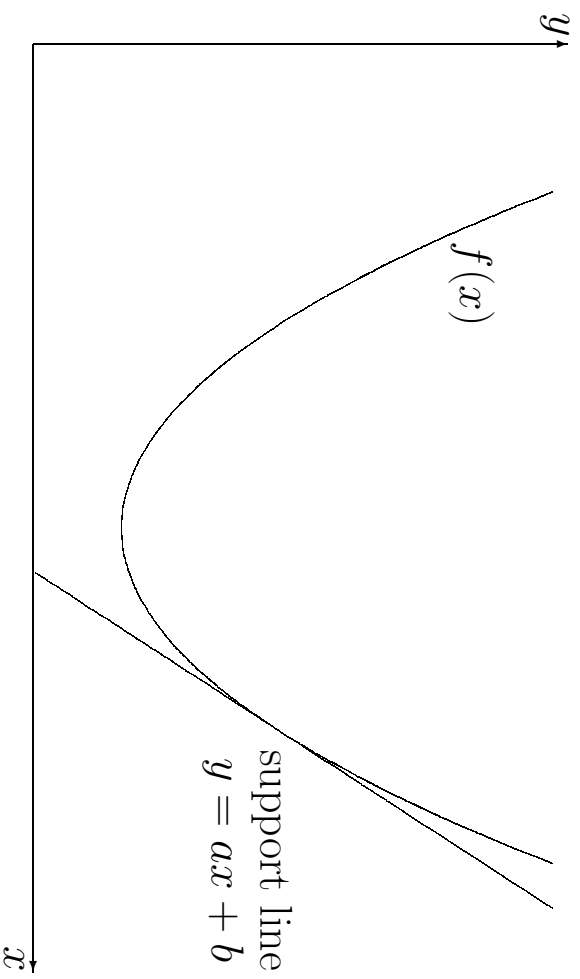
---

<sup>2</sup>A line  $y = ax + b$  is said to be a support line of function  $f(x)$  if among all lines of the same slope  $a$ , it is the largest one satisfying  $ax + b \leq f(x)$  for every  $x$ . Hence, a support line may not necessarily pass the point  $(x_0, f(x_0))$  for every  $x_0$ . Here, since we only consider convex functions, the validity of the support line at  $x_0$  passing  $(x_0, f(x_0))$  is therefore guaranteed.

## Jensen's inequality

I: b-34

The support line  $y = ax + b$  of the convex function  $f(x)$ .



# Key Notes

I: b-35

- Definitions of (weakly, strictly) stationary, ergodic and Markov (irreducible, homogeneous)
- Mode of convergences (almost surely or with probability 1, in probability, in distribution, in  $L_r$ , mean)
- Two Laws of large numbers
- Central limit theorem
- Jensen's inequality (convexity and concavity)