

# MATH5412 - Advanced Probability Theory II

Aaron Wang

[aswang@connect.ust.hk](mailto:aswang@connect.ust.hk)

Spring 2022

## Abstract

These are notes for MATH5412 at HKUST, the second course in a two part graduate-level course taught by Bao Zhigang in Spring 2022. The main focus is as a continuation of MATH5411.

## Contents

<b>1</b>	<b>February 8th, 2022</b>	<b>3</b>
1.1	Overview of the Course . . . . .	3
1.2	Heavy Tail Limiting (Poisson) Convergence . . . . .	5
1.3	Stable Law . . . . .	7
<b>2</b>	<b>February 10th, 2022</b>	<b>11</b>
2.1	Stable Law Continued . . . . .	11
2.2	Proof of Stable Law . . . . .	14
<b>3</b>	<b>February 17th, 2022</b>	<b>16</b>
3.1	Proof of Stable Law Continued . . . . .	16
3.2	Infinitely Divisible Distribution . . . . .	19
3.3	Functional Limit Theorems . . . . .	20
<b>4</b>	<b>February 22nd, 2022</b>	<b>21</b>
4.1	Functional Limit Theorems Continued . . . . .	21
4.2	Finite-Dimensional Distribution . . . . .	23
4.3	Weak Convergence in $(C, \mathcal{C})$ . . . . .	25
<b>5</b>	<b>February 24th, 2022</b>	<b>26</b>
5.1	Relative Compactness and Tightness of $\{\mu_n\}$ . . . . .	26
5.2	Examples of Convergence to Gaussian Process . . . . .	28
<b>6</b>	<b>March 1st, 2022</b>	<b>29</b>
6.1	Applications of Functional CLT . . . . .	29
<b>7</b>	<b>March 3rd, 2022</b>	<b>31</b>
7.1	Conditional Expectation . . . . .	31
7.2	Properties of Conditional Expectation . . . . .	35

<b>8</b>	<b>March 8th, 2022</b>	<b>36</b>
8.1	Properties of Conditional Expectation Cont. . . . .	36
8.2	Conditional Variance . . . . .	37
8.3	Introduction to Martingales . . . . .	37
	<b>Index</b>	<b>40</b>

# 1 February 8th, 2022

This is the first lecture of this course. We will discuss a bit about the logistics of the class and an overview of the content.

## 1.1 Overview of the Course

Although the previous course, MATH5411, covered the first half of Durrett [Dur19], this course will not be following the second half, which is largely about stochastic processes and Brownian motion, as there is another course MATH5450, Stochastic Processes, which will cover this exactly. Instead, this course will mostly look at limiting theorems, similar to the last part of MATH5411, but relaxing the i.i.d. constraints.

In MATH5411, we considered  $S_n = \sum_{i=1}^n X_i$ , but we had three assumptions:

1. The second moment  $\mathbf{E}[X^2]$  exists.
2.  $X_i$  is sequence of random variables in  $\mathbb{R}$ .
3.  $X_i$  are independent.

Now, in this course, we will attempt to relax these assumptions, and these three extensions are along completely different directions. Here is a brief overview of these three extensions:

### 1.1.1 Stable Law

From the central limit theorem, we know that if the second moment exists,  $S_n$  goes to a Gaussian distribution under an appropriate normalization. If the second moment does not exist, we have the **stable law**.

The stable law is not like Gaussian, which is universal in a sense, since as long as the second moment exists,  $S_n$  goes to a Gaussian distribution under a normalization. Once you don't have a second moment, the limiting distribution depends on the tail behavior, with different tail behaviors resulting in different limiting distributions. As such, we have a *class of distributions* when we don't have the second moment.

In the last course, besides considering  $S_n$ , we also discussed the sum of triangular array for a given sequence of random variable. In this case, the limiting distribution might not be Gaussian, with a typical example being Poisson convergence.

#### Example 1.1 (Poisson Convergence for Rare Events)

Let  $Y_{n,m} \sim \text{Be}(p_n)$  where  $p = p_n = \frac{c}{n}$  with  $Y_{n,m}$  i.i.d.  $1 \leq m \leq n$ . Then the limiting distribution of  $S_n = \sum_{m=1}^n Y_{n,m}$  approaches  $\text{Poisson}(c)$ .

If we have a sum of triangular arrays where the second moment does not exist, the possible limiting distribution is called the **infinitely divisible distribution**. It will contain the stable law as a special case.

To reiterate, in the case of triangular array, if we have the second moment, we'd get either a Gaussian or Poisson limiting distribution. For the case where we don't

have the second moment, we would get a class of distribution called the infinitely divisible distribution.

**Remark 1.2** — This part will take 3-4 lectures. References for this section can be found in Chapter 3 of Durrett [Dur19].

### 1.1.2 Functional Limiting Theorem

In the previous course, we were only concerned about the weak convergence of random variables in  $\mathbb{R}$ . What if we want to do the same for *random vectors* in  $\mathbb{R}^k$ ? Thinking even more broadly, we want to consider the weak convergence of *random functions*, or **random processes**. This leads to the second extension which is the **functional limiting theorem**.

One typical example where the functional limiting theorem is used is when considering *empirical process*.

**Example 1.3** (Example of Needing the Weak Convergence of a Random Function)

Say we have  $X \sim F$ , with  $F$  unknown and we want to perform statistical inference, with a sample  $X_1, \dots, X_n \sim F$  i.i.d. We can construct the **empirical distribution**  $F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq t)$  to approximate  $F$ .

We then want to figure out how well this approximation is by taking its difference  $F_n(t) - F(t)$ . By the law of large number, we know for any fixed  $t$ , the difference  $F_n(t) - F(t)$  goes to 0, since  $\mathbf{E}[F_n(t) - F(t)] = 0$ . We also know that the fluctuation is given by CLT if we multiply by  $\sqrt{n}$ , simply from the CLT for i.i.d. random variables.

However, we don't only want to consider this closeness for a fixed  $t$ , we want to measure closeness as a *whole function*. As such, we might introduce a distance between two functions, say the **Kolmogorov–Smirnov Statistics**  $:= \sup_t |F_n(t) - F(t)|$ . We know that this goes to zero by the *Glivenko–Cantelli theorem*<sup>a</sup>, which was introduced in the previous course. The problem is if we want to use this statistic for hypothesis testing, then we need to know the precise distribution of this statistic under suitable normalization. It turns out the suitable normalization is  $\sqrt{n}$ . If we consider  $X(t) = \sqrt{n}(F_n(t) - F(t))$ , which is a random function, the statistic becomes  $\sup_t |X(t)|$ , which is still a random variable. However, to do this, we need to find the weak limit of the whole stochastic process. Eventually,  $X(t)$  will go to the *Brownian bridge*<sup>b</sup>.

<sup>a</sup>[https://en.wikipedia.org/wiki/Glivenko%E2%80%93Cantelli\\_theorem](https://en.wikipedia.org/wiki/Glivenko%E2%80%93Cantelli_theorem)

<sup>b</sup>[https://en.wikipedia.org/wiki/Brownian\\_bridge](https://en.wikipedia.org/wiki/Brownian_bridge)

**Remark 1.4** — This part will also be quite short. References for this section can be found in Chapter 2 of Billingsley's *Convergence of Probability Measure* [Bil86].

### 1.1.3 Martingale and it's Limiting Theorem

Roughly speaking, a martingale can be thought of the sum of a random variable. This random variable, in martingale theory, are called the **martingale differences**, which are not necessarily independent. These martingale differences lie somewhere between uncorrelated and independent random variables, having more structure than uncorrelated variables, but are not as good as independent variables. As such, although they are not necessarily independent, they share many common features with independent random variables.

**Remark 1.5** — This part will be a major part of this course. References can be found in Chapter 5 of [Dur19] and Hall and Heyde's *Martingale Limit Theory and its Application* [HH80].

### 1.1.4 Concentration (if time permits)

If time permits, we will also cover something called **martingale concentration**. Very roughly speaking, concentration can be thought as an analog to the law of large numbers. Recall for WLLN, we briefly described geometric concentration. The systematic discussion of concentration will mainly focus on the non-asymptotic part, but we will still be considering a function of a large number of random variables. These random variables may be independent or not, or even martingale differences. This section is not necessarily about the limiting part of probability theory, as it focuses on the non-asymptotic behavior.

**Remark 1.6** — References for part will be taken from Vershynin's *High-Dimensional Probability* [Ver19].

## 1.2 Heavy Tail Limiting (Poisson) Convergence

Before introducing the stable law, we will quickly review the heavy tail limiting convergence from the last part of MATH5411. Heuristically, the stable law and the heavy tail convergence are very related.

As with Example 1.1, we consider a triangular array,  $Y_{n,1}, \dots, Y_{n,n} \sim \text{Be}(p)$  i.i.d. with  $p = p_n = \frac{\lambda}{n}$ . We have

$$\sum_{m=1}^n Y_{n,m} \implies \text{Poisson}(\lambda).$$

After that, we did a generalization to not require the elements in the triangular array to be i.i.d.

### Theorem 1.7 (Poisson Convergence for non i.i.d. Bernoulli Random Variables)

For each  $n$ , let  $X_{n,m}$ ,  $1 \leq m \leq n$  be independent r.v. with  $\mathbb{P}(X_{n,m} = 1) = 1 - \mathbb{P}(X_{n,m} = 0) = \beta_{n,m}$ . If  $\sum_{m=1}^n \beta_{n,m} \rightarrow \lambda$  and  $\max_m \beta_{n,m} \rightarrow 0$ , then,  $S_n = \sum_{m=1}^n X_{n,m} \implies \text{Poisson}(\lambda)$ .

**Remark 1.8** — This is similar to Lindeberg's condition for CLT.

After this, we can extend to non-Bernoulli random variables, being able to take any non-negative integer value, as long as it is “almost” Bernoulli.

**Theorem 1.9 (Poisson Convergence for non-Bernoulli Random Variables)**

For each  $n$ , let  $X_{n,m}$ ,  $1 \leq m \leq n$  be independent r.v. with  $\mathbb{P}(X_{n,m} = 1) = \beta_{n,m}$  and  $\mathbb{P}(X_{n,m} \geq 2) = \epsilon_{n,m}$ . If  $\sum_{m=1}^n \beta_{n,m} \rightarrow \lambda$ ,  $\max_m \beta_{n,m} \rightarrow 0$  and  $\sum_{m=1}^n \epsilon_{n,m} \rightarrow 0$ , then,  $S_n = \sum_{m=1}^n X_{n,m} \Rightarrow \text{Poisson}(\lambda)$ .

Now with this general result, we are able to solve a mathematical modelling problem.

**Example 1.10 (Modelling Customer Arrival)**

Suppose we open a bank and we want to know the number of arrivals  $N([s, t])$  during a time duration  $[s, t]$ . To model, this we make the following assumptions:

- (i) The number in disjoint intervals are independent
- (ii) The distribution of  $N(s, t)$  only depends on  $t - s$ , i.e. it is **time homogeneous**
- (iii)  $\mathbb{P}(N([0, h]) = 1) = \lambda h + o(h)$ , and
- (iv)  $\mathbb{P}(N([0, h]) \geq 2) = o(h)$

**Theorem 1.11**

If (i) - (iv) in Example 1.10 hold, then  $N([0, t])$  has an exact Poisson distribution with mean  $\lambda t$ .

For this example, what we really care is not the Poisson convergence, rather the consequence of this mathematical modelling problem. For this example, we not only get the distribution for a fixed  $t$ , we get a stochastic process. If we let  $t$  run from 0 to infinity, we get what is called a **Poisson point process**.

**Definition 1.12** (Poisson point process with rate  $\lambda$ ). A family of random variables  $N_t = N([0, t])$ ,  $t \geq 0$ , satisfying:

1. If  $0 = t_0 < t_1 < \dots < t_n$  then  $N_{t_k} - N_{t_{k-1}} = N([t_{k-1}, t_k])$  are all independent.
2.  $N_t - N_s \sim \text{Poisson}(\lambda(t - s))$ .

There are also a few other ways to characterize a Poisson point process, such as by the time of arrival. Thus, this process can be characterized by these points if its counting function satisfy the properties in Definition 1.12. We can also regard a Poisson point process as a random measure, leading to us being able to generalize a Poisson point process on a measure space.

**Definition 1.13** (Poisson point process on a measurable space  $(S, \mathcal{S}, \mu)$ ). A random map  $m : \mathcal{S} \rightarrow \{0, 1, \dots\}$  that for each  $\omega$  is a measure on  $\mathcal{S}$ , and has the following property:

If  $A_1, A_2, \dots, A_n$  are disjoint with  $\mu(A_i) < \infty$  then:

1.  $m(A_1), \dots, m(A_n)$  are independent.
2.  $m(A_i) \stackrel{D}{=} \text{Poisson}(\mu(A_i))$ .

where  $\mu(A) := \mathbf{E}[m(A)]$  is the mean measure of  $m$ .

### 1.3 Stable Law

Now that we have review Poisson point processes, let us move onto stable law. Consider:

$$X_1, X_2, \dots, X_n \text{ i.i.d. } S_n = \sum_{i=1}^n X_i$$

If  $\mathbf{E}X_i = \mu$  and  $\mathbf{Var}X_i = \sigma^2$ , we have:

$$\frac{S_n - n\mu}{\sqrt{n}\sigma} \Rightarrow N(0, 1)$$

Now, if  $\mathbf{E}X_i^2 = \infty$ , we want to ask if we have  $a_n, b_n, Y$  such that:

$$\frac{S_n - b_n}{a_n} \Rightarrow Y \tag{1}$$

Where  $Y$  is nondegenerate (if it is, then it would be trivial).  $a_n$  is basically the typical size of the fluctuation of  $S_n$ . In the case where the second moment exist, we know that this is of order  $\sqrt{n}$ . If we don't have the second moment, which are so called **heavy tailed random variables**, then these variables are more likely to take on large values, meaning that  $a_n$  should intuitively be larger than  $\sqrt{n}$ . How much larger depends on the explicit tail behavior of  $X_i$ . In a very special case, will be Gaussian, but in most cases it will not.

Similar to with the CLT, we eventually want to remove the assumptions about the distribution. However, let us first start with a specific special case where we know the explicit distribution of  $X_i$ . For this, we will present two solutions, the first does not have anything to do with Poisson point process, the but second will relate it to this heuristic.

Consider  $X_1, X_2, \dots$  i.i.d.

$$\mathbb{P}(X_1 > x) = \mathbb{P}(X_1 < -x) = \frac{x^{-\alpha}}{2}, \text{ for } x \geq 1, 0 < \alpha < 2.$$

The density function is thus given by:

$$f(x) = \alpha \frac{|x|^{-\alpha-1}}{2}, \quad |x| > 1$$

Note that this density function is symmetric (indicating  $b_n = 0$ ). In addition, computing the second moment using the tail sum formula, we have:

$$\mathbf{E}X_1^2 = 2 \int_1^\infty x \mathbb{P}(|X_1| > x) dx = \int_1^\infty x^{-\alpha+1} dx = \infty$$

since  $\alpha < 2$ .

**Remark 1.14** — Note that if  $\alpha = 2$ , then  $\mathbf{E}X_1^2 = 0$ . However, this is a very different case, and as such we do not consider it. In this case, the CLT holds with normalization  $\sqrt{n \log n}$ . See Theorem 1.12.3 from [UZ11].

**Remark 1.15** — As mentioned above, since this is symmetric, we can have  $b_n = 0$ . This suggests that we might have non-zero  $b_n$  for non-symmetric cases.

### Solution 1

We will try to compute limiting distribution using the Levy's continuity theorem by finding the limit of the characteristic function.

#### Theorem 1.16 (Levy's Continuity Theorem)

Suppose we have:

- a sequence of random variables  $\{X_n\}_{n=1}^\infty$ , not necessarily sharing a common probability space,
- the sequence of corresponding characteristic functions  $\{\varphi_n\}_{n=1}^\infty$ , where  $\varphi_n(t) = \mathbf{E}[e^{itX_n}]$ ,  $\forall t \in \mathbb{R}$ ,  $\forall n \in \mathbb{N}$ ,

If the sequence of characteristic functions converges pointwise to some function  $\varphi_n(t) \rightarrow \varphi(t) \forall t \in \mathbb{R}$ , then the following statements are equivalent:

- $X_n \xrightarrow{d} X$  for some random variable  $X$ .
- $\{X_n\}_{n=1}^\infty$  is tight:  $\lim_{x \rightarrow \infty} (\sup_n \mathbf{P}[|X_n| > x]) = 0$ ;
- $\varphi(t)$  is the characteristic function of some random variable  $X$ ;
- $\varphi(t)$  is a continuous function of  $t$ ;
- $\varphi(t)$  is continuous at  $t = 0$ .

We have:

$$\mathbf{E}[e^{isS_n}] = \mathbf{E}[e^{is \sum_{i=1}^n X_i}] = [\mathbf{E}[e^{isX_1}]]^n$$

Now we need to choose a normalization such that this does boil down to a characteristic function of a single point mass. In other words, we want this to be of the form  $(1 + O(\frac{1}{n}))^n$ . In this case, we choose  $\varphi(s) = \mathbf{E}[e^{isX_1}]$ , such that:

$$\mathbf{E}[e^{isS_n}] = [\mathbf{E}[e^{isX_1}]]^n = [1 - (1 - \varphi(s))]^n$$

such that  $1 - \varphi(s) \sim O(\frac{1}{n})$ . We have:

$$\begin{aligned} 1 - \varphi(s) &= \mathbf{E}[e^{isX_1}] \\ &= \int_1^\infty (1 - e^{isx}) \frac{\alpha}{2|x|^{\alpha+1}} dx + \int_{-\infty}^{-1} (1 - e^{isx}) \frac{\alpha}{2|x|^{\alpha+1}} dx \\ &= \alpha \int_1^\infty \frac{1 - \cos(sx)}{x^{\alpha+1}} dx \end{aligned} \tag{2}$$



For the case where  $s \geq 0$ , with a change of variables, we have:

$$1 - \varphi(s) = \alpha \int_1^\infty \frac{1 - \cos(sx)}{x^{\alpha+1}} dx = s^\alpha \alpha \int_s^\infty \frac{1 - \cos(u)}{u^{\alpha+1}} du \quad (3)$$

Going back to  $\frac{S_n}{a_n} = Y$ , with  $a_n$  roughly greater than  $\sqrt{n}$ . Now we absorb  $a_n$  into  $s$ , meaning that  $s$  should be really small, eventually going to zero. In the integral on the RHS of Equation 3, the singularity at infinity can be ignored, due to the  $u^{\alpha+1}$  in the denominator. For the singularity at zero, note that  $1 - \cos(u) \sim u^2$  when  $u$  close to zero, meaning the integrand is almost  $u^{1-\alpha}$ . Since  $\alpha < 2$ , this is also integrable. Thus, as  $n$  go to infinity, the integral goes to a constant, giving us  $1 - \varphi(s) = s^\alpha C_\alpha$ . Choosing  $s = \frac{t}{n^{1/\alpha}}$ , with fixed  $t \in \mathbb{R}$ . This gives us:

$$\mathbf{E} \left[ \exp \left\{ it \frac{S_n}{n^{1/\alpha}} \right\} \right] = \left[ 1 - \frac{1}{n} |t|^\alpha C_\alpha \right]^n \rightarrow e^{-C_\alpha |t|^\alpha}.$$

This means that  $\frac{S_n}{n^{1/\alpha}} \rightarrow Y$  with characteristic function  $e^{-C_\alpha |t|^\alpha}$ . This is one case of a stable law, with a specific example.

Note that in this case, Gaussian is also a special case of the stable law if we choose  $\alpha = 2$ . If we choose  $\alpha = 1$ , we would get Cauchy. In the general case, we do not have an explicit formula for the distribution or density function, so we often just express in terms of the characteristic function. For some asymptotic analysis of the density function, refer to [Nol20].

As expected, this scaling is larger than  $\sqrt{n}$  since  $\alpha < 2$ . This solution is simple because we have the explicit distribution.

## Solution 2

Before presenting the solution, we will do some preliminary analysis on the solution found above.

The reason why we expect the normalization to be larger than  $\sqrt{n}$  is because of the large tails. This motivates us to look more into these tail behavior. We believe that there is major contribution to  $S_n$  from the larger parts random variables. Investigating this way will allow us to remove any explicit assumptions of the distribution besides the tail.

Starting from Equation 2, plugging in the value of  $s$ , we want to find a  $b$  such that the contribution from the small parts:

$$\int_1^{n^b} \frac{1 - \cos(\frac{t}{n^{1/\alpha}})}{x^{\alpha+1}} dx \ll O\left(\frac{1}{n}\right).$$

We can see that this occurs if  $b < 1/\alpha$ , using by using Taylor expansion. This rough analysis tells us that the contribution from parts of the distribution before  $n^{1/\alpha}$  are not important for the limiting distribution. In other words, only the parts greater than  $n^{1/\alpha}$  are relevant to our analysis. Now let us look at the behaviour when it is on the order of  $n^{1/\alpha}$ .

By definition, for any  $b > a > 0$ , if  $a^{1/\alpha} > 1$ , then we want to consider the scale:

$$\begin{aligned}\mathbb{P}(an^{1/\alpha} < X_1 < bn^{1/\alpha}) &= \mathbb{P}(X_1 > an^{1/\alpha}) - \mathbb{P}(X_1 > bn^{1/\alpha}), \\ &= \frac{1}{2}(a^{-\alpha} - b^{-\alpha}) \cdot \frac{1}{n}\end{aligned}$$

since this is the scale where it starts to contribute to the limiting distribution. Let us study the indicator function:

$$\mathbb{1}\left(\frac{X_1}{n^{1/\alpha}} \in (a, b)\right) \sim \text{Be}\left(\frac{1}{2}(a^{-\alpha} - b^{-\alpha}) \cdot \frac{1}{n}\right).$$

This parallels the Poisson convergence. If we denote the counting measure of this indicator function, we have:

$$N_n((a, b)) = \sum_{i=1}^n \mathbb{1}\left(\frac{X_i}{n^{1/\alpha}} \in (a, b)\right) \implies N(a, b) = \text{Poisson}\left(\frac{1}{2}(a^{-\alpha} - b^{-\alpha})\right).$$

Note that we are counting the number of  $X_i$  that are of the order of  $n^{1/\alpha}$ . Note that a Poisson r.v. is of order 1, meaning that number of  $X_i$  of order  $n^{1/\alpha}$  is of constant order. This tells us that the major contribution of  $S_n$  in the heavy tail case comes from a constant number of points.

**Remark 1.17** — This is in contrast to the CLT, in which if you decompose it into small and large parts, both have significant contributions.

More generally, we can define this constant measure as:

$$N_n(A) \quad \forall A \subset \mathbb{R} \setminus (-\delta, \delta), \quad \delta n^{1/\alpha} > 1$$

then:

$$\mathbb{P}\left(\frac{X_1}{n^{1/\alpha}} \in A\right) = \int_A \frac{\alpha}{2|x|^{\alpha+1}} dx \cdot \frac{1}{n}$$

If we think of  $N_n(A)$  as a random measure, we have:

$$N_n(a) \implies N(A) \sim \text{Poisson}\left(\int_A \frac{\alpha}{2|x|^{\alpha+1}} dx\right) = \text{Poisson}(\mu(A))$$

meaning that  $N$  is a Poisson point process on  $(\mathbb{R} \setminus (-\delta, \delta), \mathcal{B}(\mathbb{R} \setminus (-\delta, \delta)), \mu(\cdot))$ .

This means that  $S_n$  will converge to the sum of points in this Poisson point process. Thus, the stable law is just the distribution of points in the Poisson point process.

**Remark 1.18** — This Poisson point process is no longer homogeneous, since the measure is not Lebesgue measurable. Most of the points in  $S_n$  will go to zero, besides the finite heavy tail ones.

## 2 February 10th, 2022

### 2.1 Stable Law Continued

Recall from the last time, we concluded that most of the contribution for  $S_n$  is from the large points of scale  $O(n^{1/\alpha})$  and that this is of constant order. Let us define an index set of large points:

$$I_n(\epsilon) = \{m \leq n : |X_m| > \epsilon n^{1/\alpha}\}$$

and define the sums:

$$\hat{S}_n(\epsilon) = \sum_{m \in I_n(\epsilon)} X_m = \sum_{m=1}^n X_m \mathbb{1}(|x_m| > \epsilon n^{1/\alpha})$$

$$\bar{S}_n(\epsilon) = S_n - \hat{S}_n(\epsilon) = \sum_{m=1}^n X_m \mathbb{1}(|X_m| \leq \epsilon n^{1/\alpha})$$

Intuitively speaking  $\hat{S}_n(\epsilon)$  represents the sum of large points and  $\bar{S}_n(\epsilon)$  represents the sum of small points.

**Remark 2.1** — Later on,  $\epsilon$  will be chosen to be as small as possible. Later we will let it go to zero along with  $n$ , e.g.  $1/\log n$ , since we might exclude relevant points. For now we will consider it fixed.

Now we have two tasks, to show

1. Show  $\frac{\bar{S}_n(\epsilon)}{n^{1/\alpha}}$  is small if  $\epsilon$  is small.
2. Find the limit of  $\frac{\hat{S}_n(\epsilon)}{n^{1/\alpha}}$ .

*Proof.* of 1.

$$\begin{aligned} \mathbf{E} \left[ \left( \frac{\bar{S}_n(\epsilon)}{n^{1/\alpha}} \right)^2 \right] &= n^{-\frac{2}{\alpha}} \cdot n \cdot \mathbf{E} \left[ (\bar{X}_1(\epsilon))^2 \right], \quad \bar{X}_i(\epsilon) = X_i \mathbb{1}(|X_i| \leq \epsilon n^{1/\alpha}) \\ \mathbf{E} \left[ (\bar{X}_1(\epsilon))^2 \right] &= \int_0^\infty 2y \mathbb{P}(|\bar{X}_1(\epsilon)| > y) dy \\ &\leq \int_0^{\epsilon n^{1/\alpha}} 2y \mathbb{P}(|X_1| > y) dy \\ &= \int_0^1 2y dy + \int_1^{\epsilon n^{1/\alpha}} 2y y^{-\alpha} dy \leq \frac{2\epsilon^{2-\alpha}}{2-\alpha} n^{\frac{2}{\alpha}-1} \end{aligned}$$

This gives us:

$$\mathbf{E} \left[ \left( \frac{\bar{S}_n(\epsilon)}{n^{1/\alpha}} \right)^2 \right] \leq \frac{2\epsilon^{2-\alpha}}{2-\alpha}, \quad 0 < \alpha < 2$$

Later we choose  $\epsilon = \epsilon_n \downarrow 0$  as  $n \rightarrow \infty$ . □

*Proof.* of 2.

Note that  $\hat{S}_n(\epsilon)$  is a sum of a random number of r.v. We will find the characteristic function using the total law of expectation:

$$\mathbf{E} \left[ \exp \left( it \frac{\hat{S}_n(\epsilon)}{n^{1/\alpha}} \right) \right] = \sum_{m=0}^n \mathbf{E} \left[ \exp \left( it \frac{\hat{S}_n(\epsilon)}{n^{1/\alpha}} \right) \middle| |I_n(\epsilon)| = m \right] \mathbb{P}(|I_n(\epsilon)| = m)$$

Now, we need to find these two terms. We will start with finding  $\mathbb{P}(|I_n(\epsilon)| = m)$ . We will use two facts:

1.  $|I_n(\epsilon)| = \sum_{m=1}^n \mathbb{1}(|X_m| > \epsilon n^{1/\alpha})$  is  $\text{Bin} \left( n, \frac{\epsilon^{-\alpha}}{n} \right) \sim \text{Poisson}(\epsilon^{-\alpha})$ , giving us  $\mathbb{P}(|X_n| > \epsilon n^{1/\alpha}) = \epsilon^{-\alpha} \frac{1}{n}$ .
2. The conditional distribution of  $\hat{S}_n(\epsilon) \middle| |I_n(\epsilon)| = m$  equals the distribution of the sum of  $m$  i.i.d. r.v. with c.d.f.  $F_\epsilon$  defined as:

$$1 - F_\epsilon(x) = \mathbb{P} \left( \frac{X_1}{n^{1/\alpha}} > x \middle| \frac{|X_1|}{n^{1/\alpha}} > \epsilon \right).$$

i.e.  $F_\epsilon$  is the conditional distribution of  $\frac{X_1}{n^{1/\alpha}}$  given  $\frac{|X_1|}{n^{1/\alpha}} > \epsilon$ .

*Proof.*

$$\begin{aligned} \mathbb{P}(\hat{S}_n(\epsilon) \in B \mid |I_n(\epsilon)| = m) &= \frac{\mathbb{P}(\hat{S}_n(\epsilon) \in B, |I_n(\epsilon)| = m)}{\mathbb{P}(|I_n(\epsilon)| = m)} \\ &= \frac{\binom{n}{m} \mathbb{P} \left( \sum_{i=1}^m X_i \in B, |X_1| > \epsilon n^{1/\alpha}, \dots, |X_m| > \epsilon n^{1/\alpha} \right)}{\binom{n}{m} \mathbb{P}(|X_1| > \epsilon n^{1/\alpha}, \dots, |X_m| > \epsilon n^{1/\alpha})} \end{aligned}$$

□

For our distribution, we have:

$$1 - F_\epsilon(x) = \frac{x^{-\alpha}}{2\epsilon^{-\alpha}}, \quad x \geq \epsilon.$$

i.e.  $F_\epsilon$  is the c.d.f. of  $\epsilon X_1$ , meaning that the characteristic function of  $F_\epsilon$  is  $\varphi(\epsilon t)$ . Consequently:

$$\begin{aligned} \mathbf{E} \left[ \exp \left\{ it \hat{S}_n(n^{1/\alpha}) \right\} \right] &= \sum_{m=0}^n \binom{n}{m} \left( \frac{\epsilon^{-\alpha}}{n} \right)^m \left( 1 - \frac{\epsilon^{-\alpha}}{n} \right)^{n-m} [\varphi(\epsilon t)]^m \\ &\rightarrow \sum_{m=0}^{\infty} \exp(-\epsilon^{-\alpha}) \cdot (-\epsilon^{-\alpha})^m \frac{[\varphi(\epsilon t)]^m}{m!} = \exp \left\{ -\epsilon^{-\alpha} (1 - \varphi(\epsilon t)) \right\} \end{aligned}$$

using Poisson approximation for binomial and DCT.

Recall earlier that we have an approximation for  $1 - \varphi(\epsilon t) = C_\alpha \epsilon^\alpha |t|^\alpha$  if  $\epsilon \rightarrow 0$ , giving us:

$$\mathbf{E} \left[ \exp \left\{ it \hat{S}_n(n^{1/\alpha}) \right\} \right] = \exp(-C_\alpha |t|^\alpha),$$

which is the same as Solution 1. Note that we need to choose  $\epsilon = \epsilon_n \downarrow 0$ . For more details see Lemma 3.7.1 of [Dur19]. □

From this solution, we can see that only the tail part matters. Now, we will try to generalize this solution.

**Definition 2.2 (slowly varying function).**  $L : \mathbb{R} \rightarrow \mathbb{R}$  is a slowly varying function if it satisfies:

$$\lim_{x \rightarrow +\infty} \frac{L(tx)}{L(x)} = 1,$$

for any fixed  $t > 0$ .

**Example 2.3**

$\log x$ ,  $\log \log x$ ,  $\log \sqrt{x}$  are slowly varying functions, but any power function  $x^t$  is not.

**Theorem 2.4 (stable law)**

Suppose  $X_1, X_2 \dots$  are i.i.d. with distribution satisfying:

- (i)  $\lim_{x \rightarrow +\infty} \mathbb{P}(X_1 > x) / \mathbb{P}(|X_1| > x) = \theta \in [0, 1]$  (tails may not be significant)
- (ii)  $\mathbb{P}(|X_1| > x) = x^{-\alpha} L(x)$ ,  $\alpha < 2$ , and  $L$  is slowly varying (general total tail)

Let  $S_n = \sum_{i=1}^n X_i$ ,  $a_n = \inf\{x : \mathbb{P}(|X_1| > x) \leq \frac{1}{n}\}$ ,  $b_n = n\mathbf{E}[X_1 \mathbb{1}(|X_1| \leq a_n)]$ , then as  $n \rightarrow \infty$ :

$$\frac{S_n - b_n}{a_n} \Rightarrow Y,$$

for a non-degenerate r.v.  $Y$ .

**Remark 2.5** —  $\theta$  in Theorem 2.4 indicates the relative heaviness between the right and left tail. If  $\theta$  close to 1, the right tail is dominant, if  $\theta \approx \frac{1}{2}$  then both tails are roughly equal.

We want to choose  $a_n$  s.t.  $\mathbb{P}\left(\frac{X_1}{a_n} \in (\alpha, \beta)\right) \sim \frac{1}{n}$  since  $\frac{S_n}{a_n} = \sum_{i=1}^n \frac{X_i}{a_n}$ , and we want the number of large points to be a constant order random variable. A natural choice is  $\mathbb{P}(|X_1| \geq a_n) \sim \frac{1}{n}$ , which gives us the quantile of  $\frac{1}{n}$ , i.e.  $a_n = \inf\{x : \mathbb{P}(|X_1| > x) \leq \frac{1}{n}\}$ .

**Remark 2.6** — We could have used  $ca_n$  for any constant  $c$ . In this case, we just choose  $c = 1$ .

For choosing  $b_n$ , we can choose  $b_n = n\mathbf{E}[X_1 \mathbb{1}(|X_1| \leq ca_n)]$  for any constant  $c$  as well. This is because  $b_n = n\mathbf{E}\left[\underbrace{X_1}_{a_n} \underbrace{\mathbb{1}(|X_1| \leq ca_n)}_{\mathbb{P}(\cdot) \sim 1/n}\right] \sim a_n$ , meaning that the limit would differ by a constant factor.

**Remark 2.7** — The reason why we can truncate to of order  $a_n$  instead of something much larger say  $a_n^2$  is because with high probability there are no

such points.

## 2.2 Proof of Stable Law

**Claim 2.8.**

$$n\mathbb{P}(|X_1| > \alpha_n) \rightarrow 1, \quad n \rightarrow \infty$$

*Proof.* omitted. □

For the tail behavior, we get:

$$\begin{aligned} n\mathbb{P}(|X_1| > x\alpha_n) &\rightarrow \theta x^{-\alpha}, \quad n \rightarrow \infty, x > 0 \\ \sim n\mathbb{P}(|X_1| > \alpha_n) \cdot \theta &= n(xa_n)^\alpha L(xa_n) \cdot \theta \\ \sim n(xa_n)^\alpha L(a_n) \cdot \theta &= nx^{-\alpha} \mathbb{P}(|X_1| > a_n) \cdot \theta \sim x^{-\alpha} \cdot \theta \end{aligned}$$

meaning that a constant in front of  $a_n$  does not affect the convergence.

This also tells us that if we use compute the counting measure:

$$N_n((x, \infty)) = \sum_{m=1}^n \mathbb{1}\left(\frac{X_m}{a_n} > x\right) \implies \text{Poisson}(\theta x^{-\alpha}).$$

More generally  $N_n(A)$  converges to a Poisson point process  $N(A)$  with mean measure

$$\mathbf{E}N(A) = \mu(A) = \int_{A \cap (0, \infty)} \theta \alpha |x|^{-(\alpha+1)} dx + \int_{A \cap (-\infty, 0)} (1 - \theta) \alpha |x|^{-(\alpha+1)} dx.$$

Now we will decompose the points into large and small parts. Let us define index set:

$$I_n(\epsilon) = \{m \leq n : |X_m| > \epsilon a_n\}$$

and define the following:

$$\begin{aligned} \hat{S}_n(\epsilon) &= \sum_{m \in I_n(\epsilon)} X_m \quad (\text{sum of large points}) \\ \bar{\mu}(\epsilon) &= \mathbf{E}[X_m \mathbb{1}(|X_n| \leq \epsilon a_n)] = \mathbf{E}\bar{X}_m(\epsilon) \\ \hat{\mu}(\epsilon) &= \mathbf{E}[X_m \mathbb{1}(\epsilon a_n < |X_n| \leq a_n)] \\ \bar{S}_n(\epsilon) &= (S_n - b_n) - (\hat{S}_n(\epsilon) - n\hat{\mu}(\epsilon)) \\ &= \sum_{m=1}^n (\bar{X}_m(\epsilon) - \bar{\mu}(\epsilon)) \quad (\text{centered sum of small points}) \end{aligned}$$

**Remark 2.9** — Unlike the special case, we need to subtract by  $b_n$ , since it is no long symmetric.

**Remark 2.10** — Note from the definition of  $b_n$ , we truncate  $\hat{\mu}(\epsilon)$  instead of going to infinity.

Now we have once again have two tasks:

1. Show  $\frac{\bar{S}_n(\epsilon)}{a_n}$  is small if  $\epsilon$  is small.
2. Find the limit of  $\frac{\hat{S}_n(\epsilon) - n\hat{\mu}(\epsilon)}{a_n}$ .

*Proof.* of 1.

$$\begin{aligned}
 \mathbf{E} \left[ \left( \frac{\bar{S}_n(\epsilon)}{a_n} \right)^2 \right] &\leq n \mathbf{E} \left[ \left( \frac{\bar{X}_1(\epsilon)}{a_n} \right)^2 \right] \\
 &\leq \int_0^\epsilon 2y \mathbb{P}(|\bar{X}_1(\epsilon)| > ya_n) dy \\
 &= \underbrace{n \mathbb{P}(|X_1| > a_n)}_{\rightarrow 1} \int_0^\epsilon 2y \underbrace{\frac{\mathbb{P}(|X_1| \geq ya_n)}{\mathbb{P}(|X_1| > a_n)}}_{y^{-\alpha}} dy \\
 &\rightarrow \int_0^\epsilon 2yy^{-\alpha} dy = \frac{2}{2-\alpha} \epsilon^{2-\alpha} \rightarrow 0 \text{ if } \epsilon \rightarrow 0
 \end{aligned}$$

□

*Proof.* of 2.

Let us first consider trying to compute the characteristic function of  $\frac{S_n}{a_n}$ , since we can add the constant part later. We have the following:

- (i)  $|I_n(\epsilon)| \rightarrow \text{Poisson}(\epsilon^{-\alpha})$
- (ii) Given  $|I_n(\epsilon)| = m$ ,  $\frac{\hat{S}_n(\epsilon)}{a_n}$  has the same distribution as the sum of  $m$  i.i.d. r.v. with c.d.f.  $F_\epsilon$ , which is again the conditional distribution of  $\frac{X_1}{a_n}$  given  $\frac{|X_1|}{a_n} \geq \epsilon$ . This time, we need to distinguish the left and right tails:

$$\begin{aligned}
 1 - F_n^\epsilon(x) &= \mathbb{P} \left( \frac{X_1}{a_n} > x \mid \frac{|X_1|}{a_n} > \epsilon \right) \rightarrow \theta \frac{x^{-\alpha}}{\epsilon^{-\alpha}} \\
 F_n^\epsilon(-x) &= \mathbb{P} \left( \frac{X_1}{a_n} < -x \mid \frac{|X_1|}{a_n} > \epsilon \right) \rightarrow (1 - \theta) \frac{|x|^{-\alpha}}{\epsilon^{-\alpha}}
 \end{aligned}$$

Let  $\Psi_n^\epsilon(t) \rightarrow \Psi^\epsilon(t)$  be the c.f. of  $F_n^\epsilon$ :

$$\Psi_n^\epsilon(t) \rightarrow \Psi^\epsilon(t) = \int_\epsilon^\infty e^{itx} \theta \epsilon^\alpha \alpha x^{-(\alpha+1)} dx + \int_{-\infty}^{-\epsilon} e^{itx} (1 - \theta) \epsilon^\alpha |x|^{-(\theta+1)} dx$$

Note that these tails only hold when  $|x| > \epsilon$ , as the density would be zero otherwise. This will be continued next lecture.

□

### 3 February 17th, 2022

#### 3.1 Proof of Stable Law Continued

Let  $\Psi_n^\epsilon(t)$  be the c.f. of  $F_n^\epsilon$ . We have:

$$\Psi_n^\epsilon(t) \rightarrow \Psi^\epsilon(t) = \int_{\epsilon}^{\infty} e^{itx} \theta \epsilon^\alpha x^{-(\alpha+1)} dx + \int_{-\infty}^{-\epsilon} e^{itx} (1-\theta) \epsilon^\alpha |x|^{-(\alpha+1)} dx$$

Now, we have:

$$\begin{aligned} \mathbf{E} \exp \left( it \frac{\hat{S}_n(\epsilon)}{a_n} \right) &= \sum_m \mathbf{E} \exp \left( it \frac{\hat{S}_n(\epsilon)}{a_n} \middle| |I_n(\epsilon)| = m \right) \cdot \mathbb{P}(|I_n(\epsilon)| = m) \\ &\sim \sum_{m=0}^{\infty} [\Psi^\epsilon(t)]^m \frac{(\epsilon^{-\alpha})^m e^{-\epsilon^{-\alpha}}}{m!} \\ &= \exp(-\epsilon^{-\alpha}(1 - \Psi^\epsilon(t))) \\ &= \exp \left[ \int_{\epsilon}^{\infty} (e^{itx} - 1) \theta \alpha x^{-(\alpha+1)} dx + \int_{-\infty}^{-\epsilon} (e^{itx} - 1) (1-\theta) \alpha |x|^{-(\alpha+1)} dx \right] \end{aligned}$$

**Remark 3.1** — The approximation should be justified by DCT, since we need to justify the convergence of the total sum.

Note that for the above case,  $\epsilon$  is fixed. In the general case, we need to send  $\epsilon \downarrow 0$ . However, when  $x \rightarrow 0$ ,  $e^{itx} - 1 \sim itx$ , and  $x \cdot x^{-(\alpha+1)} = x^{-\alpha}$  is not integrable around 0 if  $\alpha \geq 1$ .

**Remark 3.2** — When  $\theta \neq \frac{1}{2}$ , this singularity appears, which does not happen when we consider the special case.

As such, we need to consider the centered sum  $\exp \left( -it \frac{n\mu(\epsilon)}{a_n} \right)$ , with:

$$\mu(\epsilon) = \mathbf{E} X_1 \mathbb{1}(\epsilon a_n < |X_1| \leq a_n).$$

As seen previously, from the assumption of the tail behavior and slowly varying



function, we have:

$$\begin{aligned}
\mathbb{P}\left(x < \frac{X_1}{a_n} \leq y\right) &= \frac{1}{n}\theta(x^{-\alpha} - y^{-\alpha}) \\
\Rightarrow n\hat{\mu}(\epsilon)a_n &\rightarrow \int_{\epsilon}^1 x\theta\alpha x^{-(\alpha+1)}dx + \int_{-1}^{-\epsilon} x(1-\theta)\alpha|x|^{-(\alpha+1)}dx \\
\Rightarrow \mathbf{E} \exp\left(it\frac{S_n(\epsilon) - n\hat{\mu}(\epsilon)}{a_n}\right) &\rightarrow \exp\left[\int_1^{\infty} (e^{itx} - 1)\theta\alpha x^{-(\alpha+1)}dx\right. \\
&\quad + \int_{\epsilon}^1 (e^{-tx} - 1 - itx)\theta\alpha x^{-(\alpha+1)}dx \\
&\quad + \int_{-1}^{-\epsilon} (e^{itx} - 1 - itx)(1-\theta)\alpha|x|^{-(\alpha+1)}dx \\
&\quad \left. + \int_{-\infty}^{-1} (e^{itx} - 1)(1-\theta)\alpha|x|^{-(\alpha+1)}dx\right]
\end{aligned}$$

which is integrable.

Simplifying, and sending  $\epsilon \downarrow 0$ , we get:

$$\exp\left[itc + \int_0^{\infty} (e^{itx} - 1 - \frac{itx}{1+x^2})\theta\alpha x^{-(\alpha+1)}dx + \int_{-\infty}^0 (e^{itx} - 1 - \frac{itx}{1+x^2})(1-\theta)\alpha|x|^{-(\alpha+1)}dx\right]. \quad (4)$$

**Definition 3.3 (stable law).** Distribution with characteristic function of the form 4.

**Remark 3.4 (Alternative representation)** —

$$\exp[itc - b|t|^{\alpha}(1 + i\kappa \operatorname{sgn}(t)w_{\alpha} \operatorname{plha}(t))]$$

with:

$$k = 2\theta - 1 \in [-1, 1], \quad w_{\alpha} \operatorname{plha}(t) = \begin{cases} \tan(\frac{\pi\alpha}{2}), & \alpha \neq 1 \\ \frac{2}{\pi} \log |t| & \alpha = 1 \end{cases}$$

for  $0 < \alpha \leq 2$ . See (Brenman. 1968, page 204-206)

### Example 3.5

If  $\alpha = 2$ , the stable law becomes Gaussian.

### Example 3.6

If  $\alpha = 1$ ,  $c = 0$ ,  $\kappa = 0$ , we get the Cauchy distribution.

**Example 3.7**

If  $\alpha = \frac{1}{2}$ ,  $c = 0$ ,  $\kappa = 1$ ,  $b = 1$ , we get density function:

$$(2\pi y^3)^{-1/2} \exp(-1/2y), \quad y \geq 0.$$

**Remark 3.8** — The density function are not known except for the above 3 cases.

**Theorem 3.9**

$Y$  is stable law  $\iff Y$  is the weak limit of  $\frac{\sum_{i=1}^n X_i - b_n}{a_n}$  for a given sequence of i.i.d.  $X_i$ 's.

**Example 3.10**

Let  $X_1, X_2, \dots$  be i.i.d. with a density function that is symmetric about 0 and continuous and positive at 0. We claim:

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{X_i} \implies \text{a Cauchy distribution } (\alpha = 1, \kappa = 0).$$

*Proof.* Consider when  $x \rightarrow \infty$ :

$$\mathbb{P}\left(\frac{1}{X_1} > x\right) = \mathbb{P}(0 \leq X_1 < x^{-1}) = \int_0^{x^{-1}} f(y) dy = \frac{f(0)}{x}$$

Similarly, for the left tail:

$$\mathbb{P}\left(\frac{1}{X_1} < -x\right) = \frac{f(0)}{x}.$$

In addition, we have  $\theta = \frac{1}{2}$  by assumption (of symmetry), giving us  $b_n = 0$ . Thus:

$$\mathbb{P}\left(\left|\frac{1}{X_1}\right| > a_n\right) = \frac{2f(0)}{a_n} = \frac{1}{n} \implies a_n = 2f(0) \cdot n$$

Thus:

$$\frac{1}{n} \sum_{i=1}^n X_i \implies \text{Cauchy.}$$

□

**Remark 3.11** — Whenever we prove with stable law, we check the tail behavior.

Note that the centralization constant is not necessary if  $\alpha < 1$ .

Consider  $X_1, X_2, \dots$  i.i.d. with exact distribution:

$$\mathbb{P}(X_1 > x) = \theta x^{-\alpha} \quad \mathbb{P}(X_1 < -x) = (1 - \theta)x^{-\alpha}, \quad 0 < \alpha < 2, |x| \geq 1.$$

In this case, we know that  $a_n = n^{1/\alpha}$ . Meanwhile, we have:

$$\begin{aligned} b_n &= n\mathbf{E}X_1 \mathbb{1}(|X_n| < a_n) \\ &= n \int_1^{n^{1/\alpha}} (2\theta - 1)\alpha x^{-\alpha} dx \sim \begin{cases} cn & \alpha > 1 \\ cn \log n & \alpha = 1 \\ cn^{1/\alpha} & \alpha < 1 \end{cases}. \end{aligned}$$

Note that if  $\alpha < 1$ , we don't need to subtract by  $b_n$  to have convergence, but we will have a different limit if we do/don't.

**Remark 3.12** — If  $\alpha > 1$ , the constant  $cn$ .

## 3.2 Infinitely Divisible Distribution

As we mentioned previously, the stable law is the limit of  $\frac{\sum_{i=1}^n X_i - b_n}{a_n}$  for a given sequence of i.i.d.  $X_i$ 's.

On the other hand, the **infinitely divisible distribution** is the limit of  $\frac{\sum_{i=1}^n X_{n,i} - b_n}{a_n}$  for triangular array with i.i.d.  $X_{n,i}$ 's for each  $n$ .

### Example 3.13

Gaussian  $\in$  stable law, Poisson  $\in$  infinitely divisible law

Here we won't derive the infinitely divisible distributions, but we will state some results. If interested, consult the textbook.

### Example 3.14 (Poisson as an infinitely divisible distribution)

Poisson is the limit of triangular array of Bernoulli r.v.  $X_{n,1}, \dots, X_{n,n}$  with:

$$\mathbb{P}(X_{n,i} = 1) = 1 - \mathbb{P}(X_{n,i} = 0) = \frac{\lambda}{n}$$

Note that the c.f. of Poisson( $\lambda$ ) is  $\exp(\lambda(e^{it} - 1))$  which is not a stable law.

### Theorem 3.15 (Levy-Khinchin Theorem)

$Z$  has an infinitely divisible distribution  $\iff$  its c.f. is of the form:

$$\varphi(t) = \exp \left[ ict - \frac{\sigma^2 t^2}{2} + \int \left( e^{itx} - 1 - \frac{itx}{1+x^2} \right) \mu(dx) \right]$$

where  $\mu$  is a measure (not necessarily probability measure) with:

$$\mu(\{0\}) = 0, \quad \int \frac{x^2}{1+x^2} \mu(dx) < \infty.$$

**Example 3.16** (Examples of infinitely divisible distributions)

If we consider:

1. Gaussian,  $\mu = 0$  measure.
2. Poisson, we have:

$$c = \int \frac{x}{1+x^2} \mu(dx), \quad \sigma^2 = 0, \quad \mu(\{1\}) = \lambda \text{ (single point mass)}$$

3. all stable law:  $\sigma^2 = 0$ .
4. Compound Poisson:  
Let  $\xi_1, \xi_2, \dots$  be i.i.d. and  $N(\lambda)$  be an independent Poisson( $\lambda$ ) with c.f.:

$$\varphi(t) = \mathbf{E} \exp(it\xi_1) = \int \exp(itx) \mu_\xi(dx).$$

Let  $Z = \xi_1 + \dots + \xi_{N(\lambda)}$  is infinitely divisible:

$$\mathbf{E} \exp(itZ) = \exp(-\lambda(1 - \varphi(t))) = \exp \left[ \lambda \int (e^{itx} - 1) \right]$$

This is the end of this chapter about stable law.

### 3.3 Functional Limit Theorems

Our aim for this chapter is to study the weak convergence in the space  $C[0, 1]$ , which is the space of all continuous functions supported on  $[0, 1]$ .

**Remark 3.17** — The choice of considering on  $[0, 1]$  is just for convenience. We can do on other set as long as they are compact.

The weak convergence of a function means that as  $n \rightarrow \infty$ ,  $X_n(t) \rightarrow X(t)$ ,  $t \in [0, 1]$  in distribution. First, we will consider the weak convergence on a much simpler space, namely  $\mathbb{R}^k$ .

We denote a **random vector** as:

$$\vec{X} = (X^1, \dots, X^k) : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^k, \mathcal{B})$$

so that  $X^{-1}(B) \in \mathcal{F}, \forall B \in \mathcal{B}(\mathbb{R}^k)$ .

Consider a random vector sequence  $X_n = (X_n^1, \dots, X_n^k)$ ,  $n = 1, \dots$ , with c.d.f.  $F_n : \mathbb{R}^k \rightarrow [0, 1]$ :

$$F_n(\vec{x}) = \mathbb{P}(X_n^1 \leq x_1, \dots, X_n^k \leq x_k).$$

**Definition 3.18** (Convergence of a random vector sequence). We say that  $F_n$  converges to  $F$  weakly if  $F_n(x) \rightarrow F(x)$  at all continuity point of  $F$ , denoted by  $F_n \Rightarrow F$ . Further we say  $X_n$  converges to  $X$  weakly (in distribution) if  $F_n \Rightarrow F$ , denoted by  $X_n \Rightarrow X$ .

**Definition 3.19** (Alternative definition of  $X_n \Rightarrow X$ ). We say  $X_n \Rightarrow X$  if for any bounded continuous function:  $f : \mathbb{R}^k \rightarrow \mathbb{R}$ ,  $\mathbf{E}f(X_n) \rightarrow \mathbf{E}f(X)$ .

**Definition 3.20** (tightness). We say a sequence of probability measure  $\mu_n$  on  $(\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k))$  is **tight** if for any  $\epsilon > 0, \exists M = M_\epsilon > 0$ , s.t.:

$$\mu_n([-M, M]^k) \geq 1 - \epsilon.$$

**Definition 3.21** (**characteristic function of random vector**). The c.f. of  $X = (X^1, \dots, X^k)$  is defined as:

$$\varphi_X(t) = \mathbf{E} \exp(it \cdot X) = \mathbf{E} \exp \left( i \sum_{a=1}^k t_a X_a \right), \quad t = (t_1, \dots, t_k) \in \mathbb{R}^k$$

**Theorem 3.22** (inversion formula)

If  $A = [a_1, b_1] \times \dots \times [a_k, b_k]$  with  $\mu(\partial A) = 0$ :

$$\mu(A) = \lim_{T \rightarrow \infty} (2\pi)^k \int_{[-T, T]^k} \prod_{j=1}^k \left( \frac{e^{-is_j a_j} - e^{is_j b_j}}{is_j} \right) \varphi(s) ds$$

**Theorem 3.23** (continuity theorem)

Let  $X_n, 1 \leq n \leq \infty$  be a random vectors with c.f.  $\varphi_n$ , then:

$$X_n \implies X_\infty \iff \varphi_n(t) \rightarrow \varphi_\infty(t)$$

for any given  $t \in \mathbb{R}$ .

## 4 February 22nd, 2022

### 4.1 Functional Limit Theorems Continued

**Theorem 4.1** (**Cramer-Wold device**)

A sufficient condition for  $X_n \implies X_\infty$  is that:

$$\theta \cdot X_n \implies \theta \cdot X_\infty$$

for any given  $\theta \in \mathbb{R}^k$ .

*Proof.* According to the continuity theorem, we just need to show that  $\varphi_n(t) \rightarrow \varphi_\infty(t)$ . Thus, if we choose  $t = \theta$ , then this is true for any  $t$ .  $\square$

**Theorem 4.2 (multivariate CLT)**

Let  $X_1, X_2, \dots$  be i.i.d. random vectors in  $\mathbb{R}^k$  with

$$\mathbf{E}X_1 = \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix}, \quad \Gamma_{ij} = \mathbf{E}(X_n(i) - \mu_i)(X_n(j) - \mu_j),$$

If  $S_n = \sum_{i=1}^n X_i$ , then:

$$\frac{1}{\sqrt{n}}(S_n - n\mu) \Rightarrow \chi$$

which is the multivariate Gaussian  $N(0, \Gamma)$ .

**Definition 4.3 (multivariate Gaussian).** The j.p.d.f. is:

$$(2\pi)^{-\frac{k}{2}} \det(\Gamma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}X^T \Gamma X\right)$$

and the c.f. is:

$$\mathbf{E}e^{itX} = \exp\left(-\frac{1}{2}t^T \Gamma t\right)$$

*Proof.* WLOG, we assume  $\mu = 0$  by considering  $X'_n = X_n - \mu$ . Let  $\theta \in \mathbb{R}^k$ , meaning  $\theta \cdot X_n$  is a r.v. We have:

$$\mathbf{E}\theta \cdot X_n = 0, \quad \mathbf{E}(\theta \cdot X_n)^2 = \theta^T \Gamma \theta$$

By 1D CLT:

$$\frac{\sum_{i=1}^n \theta \cdot X_n}{\sqrt{n}} \Rightarrow \theta \cdot \chi$$

By the Cramer-Wold device, we have:

$$\frac{S_n}{\sqrt{n}} \Rightarrow \chi.$$

□

Now, we will return to a functional space. For this, let us review some common functional spaces:

**Definition 4.4** ( $C[0, 1]$  Functional Space). Similar to random variables, we define the space and metric:

**space:** all continuous functions on  $[0, 1]$ .

**metric/topology:** uniform topology  $\rho(x, y) = \sup_{0 \leq t \leq 1} |x(t) - y(t)|$ .

**Borel  $\sigma$ -field:** generated by open sets.

For notation, we have:  $C = C[0, 1]$ ,  $\mathcal{C} = \mathcal{B}(C[0, 1])$ .

**Definition 4.5** (random function in  $C$ ). Given a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , we say a map  $X : \Omega \rightarrow C$  a random function if:

$$X^{-1}(A) = \{\omega : X(\omega) \in A\} \in \mathcal{F}, \quad \forall A \in \mathcal{C}.$$

**Remark 4.6** — Some other possible notations include  $X_t(\omega)$ ,  $X(t; \omega)$ ,  $t \in [0, 1]$ .

For a fixed  $\omega$ , we call  $X(\cdot; \omega)$  the **trajectory/sample path/realization of  $X$** .

We can first think of  $X(t)$  as a measurable map:

$$(\Omega, \mathcal{F}, \mathbb{P}) \xrightarrow{X(t)} (C, \mathcal{C})$$

To consider it as a probability measure, we have the push forward/induced measure, giving us the distribution of  $X(t)$ .

**Definition 4.7** (Distribution of  $X(t)$ ). If we think the random map as the map between two probability space:

$$(\Omega, \mathcal{F}, \mathbb{P}) \xrightarrow{X(t)} (C, \mathcal{C}, \underbrace{\mathbb{P} \circ X^{-1}}_{\mu_X})$$

with:

$$\mu_X(A) = \mathbb{P} \circ X^{-1}(A) = \mathbb{P}(X^{-1}(A)), \quad A \in \mathcal{C}.$$

This defines a probability measure in the new measurable space.

This is a bit more abstract than probability measures on  $\mathbb{R}$  or  $\mathbb{R}^k$ , since for the simpler spaces, we have much simpler ways to identify them. For example, for  $\mathbb{R}$ , using the Stieltjes measure function, if we have a Stieltjes distribution, we don't need to know the value of the distribution on all Borel sets, we only need on all half-lines. For  $\mathbb{R}^k$ , we can use the joint distribution function. Naturally, we ask if we can identify the  $\mu_X$  in  $C[0, 1]$ . In other words, we don't want to find  $\mu_X$  on all  $A$ , but only on some special class of events. To do this, we will first introduce the finite-dimensional distribution on  $\mu_X$ .

## 4.2 Finite-Dimensional Distribution

The general idea is to pick a few time points on the sample path and consider their joint distributions.

**Definition 4.8** (natural projection). The natural projection  $\Pi_{t_1, \dots, t_k} : C \rightarrow \mathbb{R}^k$ :

$$\Pi_{t_1, \dots, t_k} = X(t) \rightarrow (X(t_1), \dots, X(t_k)), \quad t \in [0, 1]$$

with  $\mathbb{R}^k$  being equipped with the usual Euclidean metric.

Since this is a continuous map, it is also measurable. With this, we can induce a probability measure.

$$(C, \mathcal{C}, \mu_X) \xrightarrow{\Pi_{t_1, \dots, t_k}} (\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k), \mu_X \circ \Pi_{t_1, \dots, t_k}^{-1})$$

Meaning  $\mu_X \circ \Pi_{t_1, \dots, t_k}^{-1}$  is a finite dimension distribution, with:

$$\mu_X \circ \Pi_{t_1, \dots, t_k}^{-1}(B_1 \times B_2 \times \dots \times B_k) = \mu_X(\Pi_{t_1, \dots, t_k}^{-1}(B_1 \times B_2 \times \dots \times B_k))$$

**Remark 4.9** — Since  $\mu_X \circ \Pi_{t_1, \dots, t_k}^{-1}$  is on  $\mathbb{R}^k$ , we only need to consider it on rectangles.

As such, given  $\mu_X$ , we can get the probability distribution. However, what we're more interested in is the reverse, i.e. given  $\mu_X \circ \Pi_{t_1, \dots, t_k}^{-1}$ , fix a  $\mu_X$  on  $C$ . Another view of this finite-dimensional distribution is considering:

$$(\Omega, \mathcal{F}, \mathbb{P}) \xrightarrow{X_{t_1, \dots, t_k} = (X(t_1), \dots, X(t_k))} (\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k), \mathbb{P} \circ X_{t_1, \dots, t_k}^{-1})$$

Now, we have:

$$\begin{aligned} \mathbb{P} \circ X_{t_1, \dots, t_k}^{-1}(B_1 \times B_2 \times \dots \times B_k) &= \mathbb{P}((X(t_1), \dots, X(t_k)) \in B_1 \times \dots \times B_k) \\ &\implies \mathbb{P}\{\omega : (X(t_1; \omega) \dots X(t_k; \omega)) \in B_1 \times \dots \times B_k\} \end{aligned}$$

This gives us a way to get the finite dimensional distribution even if we don't know  $\mu_X$ . Now our goal is to choose arbitrary fixed points  $t_k$  and fix the  $\mu_X$ .

Suppose we have a collection of distributions  $\nu_{t_1, \dots, t_k}$  on  $\mathbb{R}^n$  for all  $t_1, \dots, t_k \in [0, 1]$ , and all  $k \in \mathbb{N}$ , we want to know when we can regard them as the class of finite-dimensional distribution of a measure  $\nu$  on  $(C, \mathcal{C})$  and whether they can uniquely describe  $\nu$ . Let us first find some requirements for this to be a finite-dimensional distribution:

1. If  $\nu_{t_1, \dots, t_k}$ 's are indeed finite-dimensional distributions of some  $\mathcal{B}$ , then:

$$\nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k) = \mathbb{P}((X(t_1), \dots, X(t_k)) \in B_1 \times \dots \times B_k)$$

for some random function  $X(t)$  with distribution  $\nu$ . Then:

$$\nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k) = \nu_{\Pi(t_1), \dots, \Pi(t_k)}(B_{\Pi(t_1)} \times \dots \times B_k),$$

where  $\Pi$  is any permutation of 1 to  $k$ .

2. We should also be able to add to  $k$  and maintain consistency:

$$\nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k) = \nu_{t_1, \dots, t_k, t_{k+1}}(B_1 \times \dots \times B_k \times \mathbb{R}).$$

It turns out that these two consistency conditions are enough.

#### Theorem 4.10 (Kolmogorov extension theorem)

Let  $T \in [0, 1]$  be some interval. For any finite sequence of distinct time  $t_1, \dots, t_k \in T$ . Let  $\mu_{t_1, \dots, t_k}$  be a probability measure on  $\mathbb{R}^k$ . Suppose these measures satisfy the two consistency conditions:

1.  $\nu_{\Pi(t_1), \dots, \Pi(t_k)}(B_{\Pi(t_1)} \times \dots \times B_k) = \nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k)$
2.  $\nu_{t_1, \dots, t_k, t_{k+1}}(B_1 \times \dots \times B_k \times \mathbb{R}) = \nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k)$ ,

Then there exists some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and a random function  $X : (\Omega, \mathcal{F}) \rightarrow (C, \mathcal{C})$  such that:

$$\nu_{t_1, \dots, t_k}(B_1 \times \dots \times B_k) = \mathbb{P}((X(t_1) \dots X(t_k)) \in B_1 \times \dots \times B_k)$$

and the distribution of  $X(t)$ , i.e.  $\nu$  is uniquely determined by  $\nu_{t_1, \dots, t_k}$ .

*Proof.* For proof, see Billingsley [Bil86]. In the book it is called Kolmogorov existence theorem.  $\square$

This means that it is enough to know all the finite-dimensional distributions.



### 4.3 Weak Convergence in $(C, \mathcal{C})$

let us consider random functions:  $X_1(t), \dots, X_n(t) \in C$  with distributions  $\mu_1, \dots, \mu_n$ .

**Definition 4.11 (weak convergence).** We say  $\mu_n \Rightarrow \mu$  or  $X_n(t) \Rightarrow X(t)$  for some  $X(t)$  with distribution  $\mu$  if:

$$\mathbf{E}f(X_n) \rightarrow \mathbf{E}f(X)$$

for any bounded and continuous functional. Here bounded means  $\sup_{X \in C} |f(X)| \leq M$  for some  $M$ .

#### Theorem 4.12 (continuous mapping theorem)

If  $h : C \rightarrow \mathbb{R}$  is continuous. Then  $X_n \Rightarrow X$  implies  $h(X_n) \Rightarrow h(X)$ , with  $X_n$  being a random function.

*Proof.* Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be any bounded continuous function. Then,  $g \circ h$  is again a bounded continuous functional. If  $X_n \Rightarrow X$ , then:

$$\mathbf{E}g(h(X_n)) \Rightarrow \mathbf{E}g(h(X))$$

by definition of. Thus,  $h(X_n) \Rightarrow h(X)$ . □

#### Example 4.13

If  $X_n \Rightarrow X$  in  $C$ , then  $\sup_{t \in [0,1]} X_n(t) \Rightarrow \sup_{t \in [0,1]} X(t)$  using triangular inequality.

#### Example 4.14

If  $X_n \Rightarrow X$  in  $C$ , then  $\int_0^1 X_n(t) dt \Rightarrow \int_0^1 X(t) dx$ .

Now, our task is to prove weak convergence in  $C$ . We first need prove the finite-dimensional convergence and tightness.

**Definition 4.15 (finite-dimensional convergence).** We say that the finite-dimensional convergence holds for  $X_n \Rightarrow X$  if for any given  $t_1, t_2, \dots, t_k$  and any  $k \in \mathbb{N}$ , there is:

$$(X_n(t_1), \dots, X_n(t_k)) \Rightarrow (X(t_1), \dots, X(t_k)).$$

**Remark 4.16 —** Finite-dimensional convergence is not convergence-determining in  $C$ . For example, let us consider:

$$Z_n(t) = nt\mathbb{1}(t \in [0, 1/n]) + (2 - nt)\mathbb{1}(t \in [1/n, 2/n])$$

we define random function:

$$X_n(t; \omega) = Z_n(t), \quad \forall \omega \in \Omega$$

hence,  $\mathbb{P}(X_n(t) = Z_n(t)) = 1$ . Let  $X(t; \omega) = 0, \forall \omega, \forall t$ . For any given

$t_1, \dots, t_k \in [0, 1]$ :

$$(X_n(t_1), \dots, X_n(t_k)) \implies (X(t_1), \dots, X(t_k))$$

However,  $X_n(t) \not\Rightarrow X(t)$ , otherwise:

$$\sup_{t \in [0,1]} X_n(t) = 1 \implies \sup_{t \in [0,1]} X(t) = 0$$

by continuous mapping theorem.

Note that for any fixed  $n$ , the finite-dimensional distribution of  $(X_n(t_1), \dots, X_n(t_k))$ 's can be used to determine the distribution of  $X_n(t)$ . However, here  $t_1, \dots, t_k$  can also be in the interval  $(0, 2/n]$ . But in finite-dimensional convergence,  $t_1, \dots, t_k$  are given at the beginning, i.e. independent of  $n$ , meaning it can jump out of the interval.

## 5 February 24th, 2022

### 5.1 Relative Compactness and Tightness of $\{\mu_n\}$

**Definition 5.1 (relative compactness).** We say  $\{\mu_n\}$  is relatively compact if for any subsequence of  $\mu_n$ , say  $\mu_{n_k}$ , one can find a further subsequence  $\mu_{n_{k_i}}$  which converges weakly.

#### Proposition 5.2

If  $\mu_n \implies \mu$  in the finite-dimensional sense and if  $\{\mu_n\}$  is relatively compact, then  $\mu_n \implies \mu$  ( $\mu_n \rightarrow \mu$  weakly).

*Proof.* Finite-dimensional convergence means that  $\mu_n \circ \Pi_{t_1, \dots, t_k}^{-1} \implies \mu \circ \Pi_{t_1, \dots, t_k}^{-1}$  for any  $k \in \mathbb{N}$  and any given  $t_1, \dots, t_k \in [0, 1]$ . Relative compactness means that for any  $\{\mu_{n_k}\}_k$ , we have  $\mu_{n_{k_i}} \implies \nu_{k_i}$  which is a probability measure. But we already know that:

$$\mu_{n_{k_i}} \circ \Pi_{t_1, \dots, t_n}^{-1} \implies \mu \circ \Pi_{t_1, \dots, t_k}^{-1}$$

On the other hand, the finite-dimensional distribution determines the measure, meaning that:  $\nu_{k_i} = \mu$ . This further means that for any subsequence  $\{\mu_{n_k}\}_k$  there exists a further  $\{\mu_{n_{k_i}}\}_i$ , s.t.  $\mu_{n_{k_i}} \implies \mu$ ,  $i \rightarrow \infty$ . meaning that  $\mu_n \implies \mu$ .  $\square$

**Definition 5.3 (tightness).** We say  $\{\mu_n\}$  is tight if for any  $\epsilon > 0$ , there exists a compact set  $K = K_\epsilon \subset C$ , s.t.:

$$\mu_n(K) \geq 1 - \epsilon, \quad n \geq n_0$$

**Remark 5.4** — Roughly speaking  $\{\mu_m\}_n$  are almost supported on a compact set of  $C$ .

#### Theorem 5.5 (Prokhorov's Theorem)

tightness  $\iff$  relative compactness.

The heuristic is that the space of a measure on a compact space is also compact. The proof is very lengthy and won't be introduced in this course.

**Remark 5.6** — If we are considering  $\mathbb{R}$  instead of  $C$ , the relation between relative compactness and tightness is easy to understand. Tightness means that it is supported by a bounded and closed set. For example if we consider  $\mu_n = \frac{1}{3}\delta_0 + \frac{2}{3}\delta_n$ . This is not tight, since the mass will escape to  $\infty$ , it is also not relatively compact, since its limit does not go to a probability measure.

Now, we will see how to show tightness. We will first use the

**Theorem 5.7 (Arzela-Ascoli Theorem)**

We consider a set  $A \subset C[0, 1]$  is relatively compact if and only if:

1.  $\sup_{x \in A} |x(0)| < \infty$  (uniform boundedness)
2.  $\lim_{\delta \rightarrow 0} \sup_{x \in A} \omega_x(\delta) = 0$ , where  $\omega_x(\delta) = \sup_{|s-t| \leq \delta} |x(s) - x(t)|$ ,  $0 < \delta \leq 1$  which is known as the **modulus of continuity**. (uniform equicontinuous)

**Definition 5.8** (**equicontinuous** at a point  $t_0$ ). For all  $x \in A, \forall \epsilon > 0, \exists \delta$  s.t.:

$$\sup_{x \in A} |x(t) - x(t_0)| \leq \epsilon \quad \forall |t - t_0| \leq \delta$$

is called equicontinuous.

**Definition 5.9** (**uniformly equicontinuous**).  $\forall \epsilon > 0, \exists \delta$  s.t.:

$$\sup_{x \in A} \sup_{|t-s| \leq \delta} |x(t) - x(s)| \leq \epsilon.$$

**Theorem 5.10**

A sequence of probability measures  $\{\mu_n\}$  on  $C$  is tight if and only if:

1. For any  $\eta > 0$ , there exists an  $a$  and  $n_0$  s.t.:

$$\mu_n\{x : |x(0)| \leq a\} \geq 1 - \eta, \quad n \geq n_0.$$

2. For each  $\epsilon > 0$  and  $\eta > 0$ , there exists a  $\delta$  and  $n_0$  s.t. :

$$\mu_n(\{x : \omega_x(\delta) \leq \epsilon\}) \geq 1 - \eta, \quad n \geq n_0.$$

We can take the intersection and closure of 1 and 2 to get a compact set by Arzela-Ascoli theorem. Proof is omitted.

We can write 1 and 2 as:

$$\mathbb{P}(|X_n(0)| > a) \leq \eta \text{ and } \mathbb{P}(\omega_{X_n}(\delta) > \epsilon) \leq \eta.$$

and this is how we can typically prove tightness, for example using Markov inequality. However, we will stop the discussion of tightness here.

## 5.2 Examples of Convergence to Gaussian Process

We will consider two examples, the first of the weak convergence of random walk to Brownian motion, and then briefly introduce a second example of the convergence of empirical processes.

**Definition 5.11 (Brownian Motion on  $[0, 1]$ ).** A 1D Brownian motion is a real valued random function  $X(t), t \in [0, 1]$  s.t.  $X(0) = 0$  and:

1. If  $0 = t_0 < t_1 < \dots < t_k$ , then  $X(t_1) - X(t_0), \dots, X(t_k) - X(t_{k-1})$  are independent.
2. If  $s, t \geq 0$ , then:  $X(s+t) - X(s) \sim N(0, t)$
3. With probability 1,  $X(t)$  is continuous.

Note that 1 and 2 give you the finite-dimensional distribution. Fixing on a rectangle, we have:

$$\begin{aligned} \mathbb{P}((X(t_1), \dots, X(t_k)) \in A_1 \times \dots \times A_k) &= \mu_x \circ \Pi_{t_1, \dots, t_k}^{-1}(A_1 \times \dots \times A_k) \\ &= \int_{A_1} dx_1 \dots \int_{A_k} dx_k \prod_{m=1}^k \beta_{t_m - t_{m-1}}(x_{m-1}, x_m), \end{aligned}$$

where:

$$\beta_t(a, b) = \frac{1}{\sqrt{2\pi t}} \exp\left(-\frac{(b-a)^2}{2t}\right)$$

### Example 5.12

If  $k = 2$ , then:

$$f_{X(t_1), X(t_2)}(x_1, x_2) \propto \exp\left(-\frac{1}{2} \begin{bmatrix} x_1 & x_2 \end{bmatrix} \Gamma^{-1} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right)$$

$$\text{with } \Gamma = \begin{bmatrix} t_1 & t_1 \\ t_1 & t_2 \end{bmatrix} \quad t_1 < t_2.$$

**Definition 5.13 (Gaussian process).** A random function whose finite-dimensional distributions are all multivariate Gaussian.

Note that Brownian motion is a Gaussian process with covariance. Also, once we identify the covariance, the process is determined.

### Theorem 5.14 (Donsker's invariance principle (functional CLT))

Let  $\xi, \dots, \xi_n$  be i.i.d. with  $\mathbf{E}\xi_i = 0$ ,  $\mathbf{Var}\xi_i = 1$ . Let  $S_n = \sum_{i=1}^n \xi_i$ ,  $S_0 = 0$ . Define a random function in  $C[0, 1]$ :

$$X_n(t) = \frac{1}{\sqrt{n}} S_{[nt]} + (nt - [nt]) \frac{1}{\sqrt{n}} \xi_{n+1}$$

note that this is rescaling the trajectory. Under the above assumption:

$$X_n(t) \implies X(t) \leftarrow \text{Brownian motion.}$$

*Proof.* We need to prove the finite-dimensional convergence and tightness. Note that tightness will be omitted, since it is case by case basis. For the finite-dim convergence, let us denote  $\varphi_{nt} = (nt - \lfloor nt \rfloor)\xi_{\lfloor nt \rfloor + 1} / \frac{1}{\sqrt{n}} \rightarrow 0$ . We have:

$$\begin{aligned} (X_n(s), X_n(t) - X_n(s)) &= \frac{1}{\sqrt{n}}(S_{\lfloor ns \rfloor, S_{\lfloor nt \rfloor} - S_{\lfloor ns \rfloor}}) + (\varphi_{ns}, \varphi_{nt} - \varphi_{ns}) \\ &\implies \left( \underbrace{N_1}_{N(0,s)}, \underbrace{N_2}_{N(0,t-s)} \right) \\ &\implies (X_n(s), X_n(t)) \implies (N_1, N_1 + N_2) \end{aligned}$$

The extension to  $k$  components is straightforward.  $\square$

## 6 March 1st, 2022

### 6.1 Applications of Functional CLT

To show an application, we ask what is the distribution of  $\max_{0 \leq i \leq n} S_i = M_n$ . To be more precise, we are asking this for a general  $\xi_i$ 's. Notice that

$$\frac{M_n}{\sqrt{n}} = \sup_{t \in [0,1]} X_n(t) \implies \sup_{t \in [0,1]} X(t)$$

by the continuous mapping theorem. Since this is a general result, we only need to find the limiting distribution of a special case, in this case the random walk. Thus, we want to know what's the maximum point the random walk has reached. Let  $\xi_i$ 's i.i.d  $\pm 1 \text{Ber}(\frac{1}{2})$ .

**Claim 6.1.**

$$\mathbb{P}(M_n \geq a) = 2\mathbb{P}(S_n > a) + \mathbb{P}(S_n = a), \quad \forall a \geq 0$$

*Proof.* Note we don't need to consider for  $a < 0$ , which is trivial, since  $S_0 = 0$ . If  $a = 0$ , we have:

$$\mathbb{P}(M_n > 0) = 1 \quad \text{since } S_0 = 0.$$

By law of total probability, we have:

$$2\mathbb{P}(S_n > 0) + \mathbb{P}(S_n = 0) = \mathbb{P}(S_n > 0) + \mathbb{P}(S_n < 0) + \mathbb{P}(S_n = 0) = 1.$$

If  $a > 0$ , we have:

$$\mathbb{P}(M_n \geq a) = \underbrace{\mathbb{P}(M_n \geq a, S_n > a)}_{\mathbb{P}(S_n > a)} + \mathbb{P}(M_n \geq a, S_n < a) + \underbrace{\mathbb{P}(M_n \geq a, S_n = a)}_{\mathbb{P}(S_n = a)}$$

What remains is to show that  $\mathbb{P}(M_n \geq a, S_n < a) = \mathbb{P}(M_n \geq a, S_n > a)$ . This is true, since there is a 1-to-1 correspondence between the paths of both sides by the reflection principle. This is simply by reflecting after the first time the trajectory reaches level  $a$ .  $\square$

Now, for any  $\alpha \geq 0$ , we set  $a_n = \lceil \alpha n^{1/2} \rceil$ . We have:

$$\begin{aligned} \mathbb{P}\left(\frac{M_n}{\sqrt{n}} \geq \alpha\right) &= \mathbb{P}(M_n \geq a_n) = 2\mathbb{P}(S_n > a_n) + \mathbb{P}(S_n = a_n) \\ &= 2\mathbb{P}(S_n/\sqrt{n} > a_n/\sqrt{n}) + \mathbb{P}(S_n/\sqrt{n} = a_n/\sqrt{n}) \\ &= 2\mathbb{P}(N(0,1) > \alpha). \end{aligned}$$

This means that the distribution function goes to:

$$\mathbb{P}\left(\frac{M_n}{\sqrt{n}} \leq \alpha\right) \rightarrow 1 - 2\mathbb{P}(N(0,1) > \alpha) = \frac{2}{\sqrt{2\pi}} \int_0^\alpha e^{-\frac{u^2}{2}} du, \quad \alpha \geq 0.$$

**Remark 6.2** — In probability theory, the general framework is to find the limiting distribution of a specific case and also prove the universal result.

Now we will consider the limiting distribution of the empirical process. This time, it converges to a Brownian bridge.

**Definition 6.3 (Brownian Bridge).** Let  $X(t), t \in [0, 1]$  be Brownian Motion. We say:

$$\mathring{X}(t) = X(t) - tX(1), \quad t \in [0, 1],$$

is a Brownian Bridge.

**Remark 6.4** — We call it a Brownian Bridge since  $\mathring{X}(0) = \mathring{X}(1) = 0$ .

**Definition 6.5** (alternate definition of Brownian Bridge). Brownian Bridge  $\mathring{X}(t)$  is a Gaussian process with:

$$\mathbf{E}(\mathring{X}(t)) = 0, \quad \mathbf{E}(\mathring{X}(s)\mathring{X}(t)) = s \wedge t - \text{s.t. } (s(1-t) \text{ if } s \leq t)$$

**Definition 6.6 (empirical distribution).** Given a r.v.  $\xi$  with c.d.f.  $F$  (often unknown in reality), we want to estimate  $F$ . Let  $\xi_i$  be i.i.d. samples of  $\xi$ . The empirical distribution is:

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\xi_i \in [0, t])$$

The empirical distribution is often used to estimate the underlying distribution, since by the SLLN, for any fixed  $t \in [0, 1]$ :

$$F_n(t) \xrightarrow{\text{a.s.}} \mathbf{E}\mathbb{1}(\xi_i \in [0, 1]) = F(t)$$

meaning it approximates it pointwise. Using the CLT, for any fixed  $t \in [0, 1]$ :

$$\sqrt{n}[F_n(t) - F(t)] \implies N(0, F(t)(1 - F(t))).$$

What we really want to know is if it can approximate it as a whole, not just pointwise convergence. This makes use of the following theorem:

**Theorem 6.7 (Glivenko-Cantelli Theorem)**

$$\|F_n - F\|_\infty = \sup_{0 \leq t \leq 1} |F_n(t) - F(t)| \xrightarrow{\text{a.s.}} 0$$

Again, this distance is called the Kolmogorov-Smirnov statistic. We might now ask what is the limiting distribution of this distance, i.e. the limiting distribution of:

$$\sup_{0 \leq t \leq 1} |\sqrt{n}(F_n(t) - F(t))|.$$

Let us define this as a random function  $X_n(t) = \sqrt{n}(F_n(t) - F(t)), t \in [0, 1]$ .

**Remark 6.8** — Note that  $X_n(t)$  is not in  $C[0, 1]$ . Instead,  $X_n(t) \in D[0, 1]$ , which is the space of functions which are right continuous with left limits. We will ignore this issue. For more rigor, check the textbook.

This is a Gaussian process, and when we choose  $t = 0$ , then  $X_n(1) = 0$ . At point  $t = 1$ , we have  $X_n(1) = 0$ . Thus, we might suspect that this is a Brownian Bridge.

### Theorem 6.9

The empirical process  $\sqrt{n}(F_n(t) - F(t))$  converges weakly to a Gaussian process  $Y(t)$  with mean 0 and covariance:

$$\mathbf{E}(Y(S)Y(t)) = F(s \wedge t) - F(s)F(t)$$

Hence,  $Y(t) \stackrel{d}{=} \dot{X}(F(t))$ . In particular, if  $F(t) = t$ , i.e.  $\xi_i$  are uniformly distributed, then  $Y(t) = \dot{X}(t)$ . Further, by continuous mapping, we have:

$$\sqrt{n} \sup_{0 \leq t \leq 1} |\sqrt{n}(F_n(t) - F(t))| \implies \sup_{0 \leq t \leq 1} \dot{X}(F(t)) = \sup_{0 \leq t < 1} \dot{X}(t)$$

if  $F$  is continuous.

This means that whenever  $F$  is continuous, the K-S statistic is non-parametric, as it does not depend on  $F$ .

**Remark 6.10** — There are other statistics that can be used, such as the Cramer Von-Mises statistics defined as:

$$\int |\sqrt{n}(F_n(t) - F(t))|^2 dF(t).$$

This is the end of this chapter on functional limiting theorems.

## 7 March 3rd, 2022

For this chapter, we will start by rigorously introducing conditional expectation, and then move on to martingale theory.

### 7.1 Conditional Expectation

In modern probability, we first give the conditional expectation in a sigma field, and then note that other cases are special cases of this. Here we will go the opposite direction. The most general definition is not constructive, meaning we usually can't write out an explicit formula, but it still maintains special properties that are useful.

Recalling the basic level of conditional expectation, considering a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , with random variable  $X = (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B})$ , the conditional expectation of  $X$  given an event  $A \in \mathcal{F}$ , with  $\mathbb{P}(A) > 0$  is:

$$\mathbf{E}(X|A) = \frac{\mathbf{E}X \mathbb{1}_A}{\mathbb{P}(A)} = \frac{\int_A X d\mathbb{P}}{\mathbb{P}(A)}.$$

**Remark 7.1** — We can write:

$$\mathbf{E}X = \frac{\mathbf{E}X \mathbb{1}_\Omega}{\mathbb{P}(A)}.$$

**Example 7.2**

The conditional probability of  $B$  given  $A$ , is:

$$\mathbb{P}(B|A) = \mathbf{E}(\mathbb{1}_B|A) = \frac{\mathbf{E}\mathbb{1}_B \cdot \mathbb{1}_A}{\mathbb{P}(A)} = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)}.$$

**Example 7.3**

The conditional expectation of  $X$  given  $Y = y$  (i.e.  $Y^{-1}(\{y\}) \in \mathcal{F}$ ), is:

$$\mathbf{E}[X|Y = y] = \psi(y)$$

since it is a function of  $y$ . This leads to a second level of conditional expectation.

The conditional expectation of  $X$  given  $Y$  is:

$$\mathbf{E}[X|Y] = \psi(Y)$$

which is a random variable.

**Remark 7.4** — The conditional expectation of a variable given an event is a number, the conditional expectation of a variable given a variable is a random variable.

To motivate this idea, consider the following example:



**Example 7.5**

Consider the probability space:

$$\Omega = \{1, 2, 3, 4, 5, 6\}, \quad \mathcal{F} = 2^\Omega, \quad \mathbb{P}(\{i\}) = \frac{1}{6}.$$

Let  $X(\omega) = \omega$  and  $Y(\omega) = \begin{cases} 1, & \omega \in \{2, 4, 6\} \\ 0, & \text{otherwise} \end{cases}$ .

- The expectation  $\mathbf{E}X$  is the "best guess" of  $X$  without any information. In this case, for example we might care about the mean square error

$$\operatorname{argmin}_e \mathbf{E}(X - e)^2 = \mathbf{E}X = 3.5.$$

- The conditional expectation given event  $Y = y$  ( $\mathbf{E}[X|Y = y]$ ) is the "best guess" given info  $Y = y$ . For example:

$$\mathbf{E}[X|Y = 1] = \frac{\mathbf{E}X \mathbb{1}_{Y=1}}{\mathbb{P}(Y = 1)} = 2 \cdot \frac{1}{6}(2 + 4 + 6) = 4.$$

On the other hand:

$$\mathbf{E}[X|Y = 0] = \frac{\mathbf{E}X \mathbb{1}_{Y=0}}{\mathbb{P}(Y = 0)} = 2 \cdot \frac{1}{6}(1 + 3 + 5) = 3.$$

one way of thinking is we are allowed to change our guess given new information.

- The conditional expectation given  $Y$  ( $\psi(Y)$ ) is given by:

$$\psi(Y) = 4\mathbb{1}(Y = 1) + 3\mathbb{1}(Y = 0).$$

In this case, we will still be given the information of  $Y$ , but we are not allowed to change our guess. Instead, we are allowed to have a "strategy" of future information you'll provide. In other words, this is the "best guess" based on the future information of  $Y$ .

**Example 7.6**

We have  $\mathbf{E}[X|X] = X$ , since the best guess based on getting the value of  $X$  is itself. On the other hand,  $\mathbf{E}[X|Y] = \mathbf{E}X$  if  $X \perp Y$ , since the future information is useless.

Going beyond these two levels, we will define the conditional expectation given a  $\sigma$ -field  $\mathcal{H} \subseteq \mathcal{F}$ . If we go back to how we explained the first two examples, the "information" generated by  $Y$  is  $\sigma(Y)$ . In the previous example, this would be  $\{\emptyset, \Omega, \{2, 4, 6\}, \{1, 3, 5\}\}$ . If  $\mathcal{H} = \sigma(Y)$ , then we shall have:

$$\mathbf{E}[X|\sigma(Y)] = \mathbf{E}[X|Y].$$

We will define this by looking at the properties of the above definition.

- The first property we see is that  $\mathbf{E}[X|Y]$  is  $\sigma(Y)$ -measurable, since the information contained cannot be more than that of  $\sigma(Y)$ .

- $\mathbf{E}[X|Y]$  is a constant on any atoms  $A \in \sigma(Y)$ , since we cannot differentiate items within the atoms.

**Definition 7.7** (atom of a *sigma*-field). There does not exist any nontrivial event  $D$  s.t.  $D \subset A$

This also means:

$$\mathbf{E}[\mathbf{E}[X|Y]]\mathbb{1}_A = \mathbf{E}[X|Y = Y(\omega)]\mathbf{E}\mathbb{1}_A = \mathbf{E}X\mathbb{1}_A,$$

Since this is true for the atoms, by linearity of expectation, it is true for any event:

$$\mathbf{E}[\mathbf{E}[X|Y]]\mathbb{1}_E = \mathbf{E}X\mathbb{1}_E, \quad \forall E \in \sigma(Y)$$

Replacing  $\sigma(Y)$  by a general  $\mathcal{H}$ , we have the following definition:

**Definition 7.8.** The conditional expectation of  $X$  given  $\mathcal{H} \subseteq \mathcal{F}$  is defined as any random variable  $Z$  satisfying:

1.  $Z$  is  $\mathcal{H}$ -measurable
2.  $\mathbf{E}Z\mathbb{1}_A = \mathbf{E}X\mathbb{1}_A$  for all  $A \in \mathcal{H}$ .

**Example 7.9**

Suppose  $X$  is  $\mathcal{H}$ -measurable. We have:  $\mathbf{E}[X|\mathcal{H}] = X$ .

**Example 7.10**

Suppose  $X$  is independent of  $\mathcal{H}$ , i.e.  $X^{-1}(B) \perp A$  for all  $B \in \mathcal{B}(\mathbb{R})$  and  $A \in \mathcal{H}$ . We have  $\mathbf{E}[X|\mathcal{H}] = \mathbf{E}X$ .

Thus, we arrive at the definition of conditional expectation:

**Definition 7.11.**  $\mathbf{E}[X|Y] = \mathbf{E}[X|\sigma(Y)]$ . For an event, given an event  $A \in \mathcal{F}$ , let  $\mathcal{H} = \{\emptyset, \Omega, A, A^C\}$ , we have:

$$\mathbf{E}[X|\mathcal{H}] = \frac{\mathbf{E}X\mathbb{1}_A}{\mathbb{P}(A)} \cdot \mathbb{1}_A + \frac{\mathbf{E}X\mathbb{1}_{A^C}}{\mathbb{P}(A^C)} \cdot \mathbb{1}_{A^C}$$

If we observe event  $A$ , then this reduces.

Now, we will consider the existence and uniqueness. First, for uniqueness, we will show that any two variables that satisfy these two properties are the same almost surely. Suppose we have  $Z$  and  $\tilde{Z}$  which both satisfy 1 and 2. For any  $\epsilon > 0$ , let:

$$A_\epsilon = \{\omega : Z(\omega) - \tilde{Z}(\omega) > \epsilon\} \in \mathcal{H}$$

by property 1. By property 2, we have:

$$0 = \int_{A_\epsilon} (Z - \tilde{Z}) d\mathbb{P} = \int_{A_\epsilon} (Z - \tilde{Z}) d\mathbb{P} \geq \epsilon \mathbb{P}(A_\epsilon) \implies \mathbb{P}(A_\epsilon) = 0 \implies Z = \tilde{Z} \text{ a.s.}$$

Now we will consider the existence.

**Definition 7.12.** A measure  $\nu$  is absolutely continuous w.r.t  $\mu$ , written as  $\nu \ll \mu$  if  $\mu(A) = 0 \implies \nu(A) = 0$ .

**Example 7.13**

Consider  $\mu$  be the Lebesgue measure and  $\nu$  be a distribution of a continuous r.v. Then:

$$\nu(A) = \int_A f d\mu$$

meaning that  $\nu \ll \mu$ .

**Theorem 7.14 (Radon-Nikodym Theorem)**

Let  $\mu$  and  $\nu$  be  $\sigma$ -finite measures on  $(\Omega, \mathcal{F})$ . If  $\nu \ll \mu$ , then there is a function  $f$  which is  $\mathcal{F}$ -measurable, and for all  $A \in \mathcal{F}$ :

$$\nu(A) = \int_A f d\mu,$$

with  $f$  being the **Radon-Nikodym derivative** often denoted by  $\frac{d\nu}{d\mu}$ .

Now we will prove the existence of  $\mathbf{E}[X|\mathcal{H}]$ :

*Proof.* Let the original probability space be  $(\Omega, \mathcal{F}, \mathbb{P})$ . Consider  $X \geq 0$ . Let  $\mu = \mathbb{P}$  and set:

$$\nu(A) = \int_A X d\mathbb{P} = \mathbf{E}X \mathbb{1}_A, \quad \forall A \in \mathcal{H}$$

apparently  $\nu \ll \mu$ . By R-N theorem, we know that  $\frac{d\nu}{d\mu}$  is  $\mathcal{H}$ -measurable and:

$$\nu(A) = \int_A X d\mu = \mathbf{E}X \mathbb{1}_A$$

but also:

$$\nu(A) = \int_A \frac{d\nu}{d\mu} d\mu = \mathbf{E}[\mathbf{E}[X|\mathcal{H}]] \mathbb{1}_A$$

□

## 7.2 Properties of Conditional Expectation

Some of the main properties of conditional expectation can be categorized into those related to measurability, and those related to independence.

### 7.2.1 Pulling out independent factors

- If  $X \perp \mathcal{H}$ ,  $\mathbf{E}[X|\mathcal{H}] = \mathbf{E}X$ .
- If  $X$  is independent of  $\sigma(Y, \mathcal{H})$ ,  $\mathbf{E}[XY|\mathcal{H}] = \mathbf{E}X \mathbf{E}[Y|\mathcal{H}]$ . (note needs both independent of  $Y$  and  $\mathcal{H}$ ).
- If  $X \perp Y$ ,  $\mathcal{G} \perp \mathcal{H}$ ,  $X\mathcal{H}$ ,  $X\mathcal{G}$ :

$$\mathbf{E}[\mathbf{E}[XY|\mathcal{G}]] = \mathbf{E}X \mathbf{E}Y = \mathbf{E}[\mathbf{E}[XY|\mathcal{H}]]$$

*Proof.* Note that  $\mathbf{E}[XY|\mathcal{G}]$  is  $\mathcal{G}$ -measurable, so  $\mathbf{E}[XY|\mathcal{G}] \perp \mathcal{H}$ . The proof then follow using the law of total expectation,  $\mathbf{E}[\mathbf{E}[XY|\mathcal{G}]] = \mathbf{E}[XY]$ . □

### 7.2.2 Pulling out known factors

- If  $X$  is  $\mathcal{H}$ -measurable,  $\mathbf{E}[X|\mathcal{H}] = X$ .
- $\mathbf{E}[f(X)|\mathcal{H}] = f(X)$ .
- If  $X$  is  $\mathcal{H}$ -measurable, then  $\mathbf{E}[XY|\mathcal{H}] = X\mathbf{E}[Y|\mathcal{H}]$ .
- $\mathbf{E}[f(X)Y|X] = f(X)\mathbf{E}[Y|X]$ .

## 8 March 8th, 2022

### 8.1 Properties of Conditional Expectation Cont.

**Definition 8.1** (tower property). For sub  $\sigma$ -field  $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \mathcal{F}$ , we have:

$$\mathbf{E}[\mathbf{E}[X|\mathcal{H}_2]|\mathcal{H}_1] = \mathbf{E}[\mathbf{E}[X|\mathcal{H}_1]|\mathcal{H}_2] = \mathbf{E}[X|\mathcal{H}_1]$$

*Proof.* Note that  $\mathbf{E}[\mathbf{E}[X|\mathcal{H}_1]|\mathcal{H}_2] = \mathbf{E}[X|\mathcal{H}_1]$  is trivial, since  $\mathbf{E}[X|\mathcal{H}_1]$  is  $\mathcal{H}_1$ -measurable so is  $\mathcal{H}_2$ .

Second, denote  $Y = \mathbf{E}[X|\mathcal{H}_1]$ . We have that:

$$\mathbf{E}[\mathbf{E}[X|\mathcal{H}_2]|\mathcal{H}_1] = \mathbf{E}[X|\mathcal{H}_1] \iff \mathbf{E}[Y|\mathcal{H}_1] = \mathbf{E}[X|\mathcal{H}_1].$$

By the definition of  $\mathbf{E}[Y|\mathcal{H}_1]$ , it shall be  $\mathcal{H}_1$ -measurable and  $\mathbf{E}[\mathbf{E}[Y|\mathcal{H}_1]\mathbb{1}_A] = \mathbf{E}Y\mathbb{1}_A$ ,  $A \in \mathcal{H}_1$ . Thus, we need to show that:

1.  $\mathbf{E}[X|\mathcal{H}_1]$  is  $\mathcal{H}_1$ -measurable (trivial)
2.  $\mathbf{E}[\mathbf{E}[X|\mathcal{H}_2]\mathbb{1}_A] = \mathbf{E}Y\mathbb{1}_A$  for all  $A \in \mathcal{H}_1$ . We claim that both sides are equal to  $\mathbf{E}X\mathbb{1}_A$ . The LHS is true for any  $A \in \mathcal{H}_1$  by definition, and the RHS is true for any  $A \in \mathcal{H}_2$  by definition of  $\mathbf{E}[X|\mathcal{H}_2]$ .

□

#### Example 8.2

The **law of total expectation** is a special case, as we have:

$$\mathbf{E}[\mathbf{E}[Z|\mathcal{H}]] = \mathbf{E}[\mathbf{E}[Z|\mathcal{H}|\{\emptyset, \Omega\}]] = \mathbf{E}[Z|\{\emptyset, \Omega\}] = \mathbf{E}Z.$$

Note that  $\mathbf{E}[Z|\{\emptyset, \Omega\}] = \mathbf{E}Z$ .

#### Corollary 8.3

Here are some corollaries of the tower property:

- Law of total expectation
- If  $Y$  is  $\mathcal{H}$ -measurable, then  $\mathbf{E}[\mathbf{E}[X|\mathcal{H}]] = \mathbf{E}[X|Y]$ .
- $\mathbf{E}[\mathbf{E}[X|Y]|f(Y)] = \mathbf{E}[X|f(Y)]$
- $\mathbf{E}[\mathbf{E}[X|Y, Z]|Y] = \mathbf{E}[X|Y]$

### 8.1.1 Additional Properties of Conditional Expectation

**Linearity:**  $\mathbf{E}[aX_1 + bX_2|\mathcal{H}] = a\mathbf{E}[X_1|\mathcal{H}] + b\mathbf{E}[X_2|\mathcal{H}]$

**Monotonicity:**  $\mathbf{E}[X_1|\mathcal{H}] \leq \mathbf{E}[X_2|\mathcal{H}]$  if  $X_1 \leq X_2$

**Jensen:** If  $\varphi : \mathbb{R} \rightarrow \Omega$  is convex, then  $\varphi(\mathbf{E}[X|\mathcal{H}]) \leq \mathbf{E}[\varphi(X)|\mathcal{H}]$

## 8.2 Conditional Variance

**Definition 8.4 (conditional variance).** The conditional variance is defined by:

$$\text{Var}[X|\mathcal{H}] = \mathbf{E}[(X - \mathbf{E}[X|\mathcal{H}])^2|\mathcal{H}] = \mathbf{E}[X^2|\mathcal{H}] - (\mathbf{E}[X|\mathcal{H}])^2$$

**Theorem 8.5 (law of total variance)**

$$\text{Var}(X) = \underbrace{\mathbf{E}(\text{Var}[X|\mathcal{H}])}_{\mathbf{E}(X - \mathbf{E}[X|\mathcal{H}])^2} + \underbrace{\text{Var}(\mathbf{E}[X|\mathcal{H}])}_{\mathbf{E}(\mathbf{E}[X|\mathcal{H}] - \mathbf{E}X)^2}$$

## 8.3 Introduction to Martingales

Consider a discrete stochastic process  $\{X(t), t \in T\}$ , with  $T$  being discrete.

**Definition 8.6 (filtration).** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space. A filtration  $\{\mathcal{F}_n\}_{n \geq 1}$  is a sequence of  $\sigma$ -fields, s.t.  $\mathcal{F}_i \subseteq \mathcal{F}_{i+1} \subseteq \mathcal{F}$  for all  $i$ .

Let  $X_1, X_2, \dots, X_n$  be a random process on  $(\Omega, \mathcal{F}, \mathbb{P})$ . Let  $\mathcal{F} = \sigma(X_1, \dots, X_i)$

**Definition 8.7 (adapted sequence).** A random process  $\{X_n\}$  is adapted to  $\{\mathcal{F}_n\}$  if  $X_n$  is  $\mathcal{F}_n$ -measurable for all  $n$ .

### Example 8.8

$\mathcal{F}_n$ : a collection of info of stock market up to day  $n$ .  $X_n$ : stock price of day  $n$ .

**Definition 8.9 (martingale).** If  $X_n$  is a sequence of r.v. and  $\{\mathcal{F}_n\}$  is a filtration, with:

- $\mathbf{E}|X_n| < \infty$
- $\{X_n\}$  is adapted to  $\{\mathcal{F}_n\}$
- $\mathbf{E}[X_{n+1}|\mathcal{F}_n] = X_n$  for all  $n$

then we call  $\{X_n\}$  to be a martingale w.r.t  $\{\mathcal{F}_n\}$ .

**Definition 8.10 (submartingale and supermartingale).** If the "=" is replaced with " $\geq$ ", then it is called a submartingale. If it's replaced with a " $\leq$ " it is called a supermartingale.

**Remark 8.11 —** Submartingale means that the expectation is greater than previous, meaning that you're playing a favorable game. However, most casinos are supermartingales.

**Example 8.12**

Let us consider successive tosses of a fair coin:

$$\xi_n = \begin{cases} 1 & \text{if } n\text{-th coin is H} \\ -1 & \text{if } n\text{-th coin is T} \end{cases}, \quad \mathbf{E}\xi_n = 0.$$

Let  $X_n = \sum_{i=1}^n \xi_i$ . Then  $\{X_n\}$  is a martingale w.r.t.  $\{\mathcal{F}_n\}$  with  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ .

**Example 8.13**

(**Polya's urn model**) Consider a urn with  $r$  red and  $g$  green balls. At each time we draw a ball and replace it with  $c + 1$  balls of the same color drawn. Let  $X_n$  be the fraction of green balls after the  $n$ -th draw. We claim that  $X_n$  is a martingale w.r.t.  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ .

*Proof.* We have:

$$\mathbf{E}[X_{n+1}|\mathcal{F}_n] = \mathbf{E}[X_{n+1}|X_n].$$

Suppose at step  $n$ , there are  $i$  red and  $j$  green. Then:

$$X_{n+1} = \begin{cases} \frac{j+c}{i+j+c}, & \text{with probability } \frac{j}{j+i} \\ \frac{j}{i+j+c}, & \text{with probability } \frac{i}{j+i}, \end{cases}$$

Thus:

$$\mathbf{E}\left[X_{n+1}|X_n = \frac{j}{i+j}\right] = \frac{j+c}{i+j+c} \frac{j}{i+j} + \frac{j}{i+j+c} \frac{i}{i+j} = \frac{j}{i+j}.$$

□

**Example 8.14**

(**Galton-Watson process**) Let  $\xi_i^n$ ,  $i, n \geq 1$  be i.i.d. nonnegative integer-valued. Let  $z_n$   $n > 0$  be defined as:

$$Z_0 = 1 \quad Z_{n+1} = \begin{cases} \xi_i^{n+1} + \dots + \xi_{Z_n}^{n+1} & \text{if } Z_n \neq 0 \\ 0 & \text{if } Z_n = 0 \end{cases}$$

Here  $\xi_i^{n+1}$  is the number of offspring of  $i$ -th individual in  $n$ -th generation. Let  $\mathcal{F}_n = \sigma\{\xi_i^m, i \geq 1, 1 \leq m \leq n\}$  and  $\mu : \mathbf{E}\xi_i^m \in (0, \infty)$ . Then:

$$\frac{Z_n}{\mu^n} \text{ is a martingale w.r.t. } \mathcal{F}_n.$$

*Proof.* We claim that  $\mathbf{E}Z_n = \mu^n$ , since:

$$\begin{aligned}\mathbf{E}Z_{n+1} &= \sum_{k=0}^{\infty} \mathbf{E} \left( \sum_{i=1}^k \xi_i^{n+1} | Z_n = k \right) \mathbb{P}(Z_n = k) \\ &= \sum_{k=0}^{\infty} \left( \mathbf{E} \sum_{i=1}^k \xi_i^{n+1} \right) \cdot \mathbb{P}(Z_n = k) \\ &= \sum_{k=0}^{\infty} k \mathbb{P}(Z_n = k) \cdot \mu = \mu \mathbf{E}Z_n.\end{aligned}$$

To show that it's a martingale, we have:

$$\begin{aligned}\mathbf{E}(Z_{n+1} | \mathcal{F}_n) &= \sum_{k=1}^{\infty} \mathbf{E}(Z_{n+1} \mathbb{1}(Z_n = k) | \mathcal{F}_n) \\ &= \sum_{k=1}^{\infty} \mathbf{E} \left( \sum_{i=0}^k \xi_i^{n+1} \mathbb{1}(Z_n = k) | \mathcal{F}_n \right) \\ &= \sum_{k=1}^{\infty} k \mathbb{1}(Z_n = k) \mu = \mu Z_n\end{aligned}$$

□

### Example 8.15

(**de Moivre's martingale**) Consider an unfair coin with probability  $p$  of H:

$$\xi_n = \begin{cases} 1 & n\text{-th flipping is H} \\ 0 & \text{otherwise} \end{cases}$$

Let  $X_n = \sum_{i=1}^n \xi_i$  and  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ . Let  $Y_n = \left(\frac{1-p}{p}\right)^{X_n}$ , then  $Y_n$  is a martingale w.r.t.  $\mathcal{F}_n$ .

*Proof.*

$$\mathbf{E}[Y_{n+1} | \mathcal{F}_n] = p \cdot \left(\frac{1-p}{p}\right)^{X_n+1} + (1-p) \left(\frac{p}{p}\right)^{X_n-1} = \left(\frac{1-p}{p}\right)^{X_n} = Y_n$$

□

## References

- [HH80] P. Hall and C.C. Heyde. *Martingale Limit Theory and Its Application*. Probability and mathematical statistics. Academic Press, 1980. ISBN: 9781483240244. URL: <http://www.stat.yale.edu/~mjk56/MartingaleLimitTheoryAndItsApplication.pdf>.
- [Bil86] Patrick Billingsley. *Probability and Measure*. Second. John Wiley and Sons, 1986. URL: <https://www.colorado.edu/amath/sites/default/files/attached-files/billingsley.pdf>.
- [UZ11] V.V. Uchaikin and V.M. Zolotarev. *Chance and Stability: Stable Distributions and their Applications*. Modern Probability and Statistics. De Gruyter, 2011. ISBN: 9783110935974. URL: <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.117.6505&rep=rep1&type=pdf>.
- [Dur19] Rick Durrett. *Probability: Theory and Examples*. 2019. URL: [https://services.math.duke.edu/~rtd/PTE/PTE5\\_011119.pdf](https://services.math.duke.edu/~rtd/PTE/PTE5_011119.pdf).
- [Ver19] Roman Vershynin. *High-Dimensional Probability*. 2019. URL: <https://www.math.uci.edu/~rvershyn/papers/HDP-book/HDP-book.pdf>.
- [Nol20] J.P. Nolan. *Univariate Stable Distributions: Models for Heavy Tailed Data*. Springer Series in Operations Research and Financial Engineering. Springer International Publishing, 2020. ISBN: 9783030529154. URL: <https://link.springer.com/content/pdf/10.1007/978-3-030-52915-4.pdf>.



# Index

- adapted sequence, 37
- Arzela-Ascoli Theorem, 27
- Brownian Bridge, 30
- Brownian Motion on  $[0, 1]$ , 28
- characteristic function of random vector, 21
- conditional variance, 37
- continuous mapping theorem, 25
- Cramer-Wold device, 21
- de Moivre's martingale, 39
- Donsker's invariance principle, 28
- empirical distribution, 4, 30
- equicontinuous, 27
- filtration, 37
- functional limiting theorem, 4
- Galton-Watson process, 38
- Gaussian process, 28
- Glivenko-Cantelli Theorem, 30
- heavy tailed random variables, 7
- infinitely divisible distribution, 3, 19
- Kolmogorov extension theorem, 24
- Kolmogorov-Smirnov Statistics, 4
- law of total expectation, 36
- Levy's Continuity Theorem, 8
- Levy-Khinchin Theorem, 19
- martingale, 37
- martingale concentration, 5
- martingale differences, 5
- modulus of continuity, 27
- multivariate CLT, 22
- multivariate Gaussian, 22
- Poisson point process, 6
- Polya's urn model, 38
- Possion convergence, 5
- Prokhorov's Theorem, 26
- Radon-Nikodum derivative, 35
- Radon-Nikodym Theorem, 35
- random processes, 4
- random vector, 20
- relative compactness, 26
- slowly varying function, 13
- stable law, 3, 13, 17
- submartingale, 37
- supermartingale, 37
- tight, 21
- tightness, 26
- time homogeneous, 6
- trajectory/sample path/realization of  $X$ , 23
- uniformly equicontinuous, 27
- weak convergence, 25