

STAT575
Large Sample Theory

Pingbang Hu

February 25, 2024

Abstract

This is a graduate-level theoretical statistics course taught by [Georgios Fellouris](#) at University of Illinois Urbana-Champaign, aiming to provide an introduction to asymptotic analysis of various statistical methods, including weak convergence, Lindeberg-Feller CLT, asymptotic relative efficiency, etc.

We list some references of this course, although we will not follow any particular book page by page: *Asymptotic Statistics* [[Vaa98](#)], *Asymptotic Theory of Statistics and Probability* [[Das08](#)], *A course in Large Sample Theory* [[Fer17](#)], *Approximation Theorems of Mathematical Statistics* [[Ser09](#)], and *Elements of Large-Sample Theory* [[Leh04](#)].



This course is taken in Spring 2024, and the date on the cover page is the last updated time.

Contents

| | | |
|----------|---|----------|
| 1 | Introduction | 2 |
| 1.1 | Parametrized Approach | 2 |
| 1.2 | Hypothesis Testing | 3 |
| 2 | Modes of Convergence | 4 |
| 2.1 | Different Modes of Convergence | 4 |
| 2.2 | Weak Convergence | 7 |
| 2.3 | Characteristic Function | 22 |
| 2.4 | Fundamental Theorems of Probability | 27 |

Chapter 1

Introduction

Lecture 1: Introduction to Large Sample Theory

Say we first collect n data points $x_1, \dots, x_n \in \mathbb{R}^d$, large sample theory concerns with the limiting theory as $n \rightarrow \infty$. We may treat x_i as a realization of a random vector X_i on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. In this course, we will primarily consider the case that X_i 's are i.i.d., i.e., independent and identically distributed from a distribution function, or the *cumulative density function* (cdf) F such that 16 Jan. 9:30

$$X = (X^1, \dots, X^d) \sim F(x_1, \dots, x_d) \equiv \mathbb{P}(X^1 \leq x_1, \dots, X^d \leq x_d)$$

for all $x_i \in \mathbb{R}$. If we have access to F , we can compute the corresponding *probability density function* (pdf) f , and then have access to $\mathbb{P}(X \in A)$ for all (measurable) $A \subseteq \mathbb{R}^d$ of interest.

Notation. In the measure-theoretic sense, the measure \mathbb{P} in $(\Omega, \mathcal{F}, \mathbb{P})$ is the **Lebesgue-Stieltjes measure** μ_F induced by the distribution function F . When doing integration, we will often denote

$$d\mu_F(x) = d\mathbb{P}(x) =: F(dx) =: dF(x) =: f(x)dx$$

Remark. If we know any of the above, we know every thing about the population.

Hence, the goal is to compute this by collecting data x_i 's, which is a statistical inference problem.

1.1 Parametrized Approach

There are various ways of doing this task, one way is the so-called parametrized approach. By postulating a family of cdfs $\{F_\theta, \theta \in \Theta\}$ where Θ is often a subset of \mathbb{R}^m for some m (generally $\neq n$), the goal is to select a member of this family that is the “closest”, or the “best fit” to the truth, i.e., F , based on the data.

Note. To emphasize that this depends on the data, we sometimes write the function we found as $\hat{\theta}_n(x_1, \dots, x_n)$ so that $F_{\hat{\theta}_n(x_1, \dots, x_n)}$ is our proxy for F .

Now, assume that the family is initially given, the problem is then how to select $\hat{\theta}_n$.

Example. Fisher suggested that we should look at the maximum likelihood estimator (MLE).

The justification for MLE is not about finite n , but about its asymptotic behavior when $n \rightarrow \infty$. Specifically, we have the following theorem due to Fisher (informally stated).

Theorem 1.1.1 (Fisher). If $F \in \{F_\theta: \theta \in \Theta\}$, i.e., if $F = F_{\theta^*}$ for some $\theta^* \in \Theta$, then under certain conditions, $\hat{\theta}_n$ will be “close” to θ^* as $n \rightarrow \infty$. Under some other conditions, $\sqrt{n}(\hat{\theta}_n - \theta)$ is approximately Gaussian with variance being the “best possible” in some sense.

On the other hand, in the misspecified case, i.e., $F \notin \{F_\theta, \theta \in \Theta\}$, we can still compute the MLE, which leads to another justification for MLE since even in this case, $\hat{\theta}_n$ will still be “close” to θ^* such that F_{θ^*} is, in some sense, the “closest” to F among all possible F_θ (minimizing divergence, to be precise).

1.2 Hypothesis Testing

We will also develop theory for hypothesis testing for some hypothesis we’re interested in, e.g., whether the data we collect is really i.i.d., or whether our proposed family is reasonable enough. Say now X_i ’s are scalar random variable with $\mathbb{E}[X] = \mu$, and we want to test the null hypothesis $H_0: \mu = 0$.

Example. Consider a controlled group Z and a treatment group Y , and we observe Z_1, \dots, Z_n , and Y_1, \dots, Y_n , respectively, and compute $X_i = Z_i - Y_i$ for all i . Testing H_0 on the distribution of X will show the effect of the treatment.

To do this, a well-known method is the so-called t -test. Let s_n to be the sample standard derivation, then we can compute

$$T_n = \frac{\bar{X}_n}{s_n/\sqrt{n}} \sim t_{n-1}$$

as long as X is Gaussian, i.e., the t -statistics for H_0 . What if X is not an Gaussian? We will show that even if X is not Gaussian, this result is “approximately valid” when n is “large enough” as long as $\text{Var}[X] < \infty$.

Remark (Sample Size). When we say n is “large enough”, what we mean really depends on how fast the underlying distribution will approach Gaussian as n grows. Hence, if we can say more about the underlying population, we can say more about when does n is “large enough”; otherwise such a limiting theory might be completely useless in practice.

What if now $\text{Var}[X]$ doesn’t exist? When the population has a heavy tail distribution, then second moment may not exist.

Example (Cauchy distribution). The Cauchy distribution doesn’t have finite moment of order greater than 1.

In this case, some other test is needed. A simple test would be looking at the sign of X_i , i.e., the sign test.

Example (Sign test). We might reject H_0 if $\sum_{i=1}^n \mathbb{1}_{X_i > 0}$ is large. Note that under H_0 , $\sum_{i=1}^n \mathbb{1}_{X_i > 0} \sim \text{Bin}(n, 1/2)$, and this test is valid even if expectation doesn’t exist.

We see that without saying anything about F , the sign test is valid even for $n = 3$ or 5 as the sum is exactly binomial distribution under H_0 . Although simple and have good property, only looking at the sign of X_i might be too weak. A natural idea is to look at the absolute value of X_i .

Example (Wilcoxon’s rank-sum test). Let $R_{i,n}$ to be the rank of $|X_i|$, then consider the so-called *Wilcoxon’s rank-sum test*

$$\sum_{i=1}^n \mathbb{1}_{X_i > 0} R_{i,n}.$$

As one can imagine, the closed form of the above sum will be complicated; however, asymptotically, the above statics will follow Gaussian again, such that the rate of convergence doesn’t depend on the underlying population.

Finally, we also ask how can we compare these different tests? This will also be addressed in this course.

Chapter 2

Modes of Convergence

Lecture 2: Modes of Convergence

2.1 Different Modes of Convergence

18 Jan. 9:30

Given a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, consider a sequence of d -dimensional random vectors (X_n) and a random vector X , i.e., $X_n, X: \Omega \rightarrow \mathbb{R}^d$. We now discuss different modes of convergence for (X_n) .

Definition 2.1.1 (Point-wise converge). (X_n) *point-wise converges* to X , denoted as $X_n \rightarrow X$, if $X_n(\omega) \rightarrow X(\omega)$ for all $\omega \in \Omega$.^a

^aI.e., for every $\epsilon > 0$, there exists $n_0(\omega) \in \mathbb{N}$ such that for every $n \geq n_0$, $\|X_n(\omega) - X(\omega)\|_2 < \epsilon$.

Since we don't care about measure zero sets, we may instead consider the following.

Definition 2.1.2 (Converge almost-surely). (X_n) *converges almost-surely* to X , denoted as $X_n \xrightarrow{\text{a.s.}} X$, if $\mathbb{P}(X_n \rightarrow X) = 1$.^a

^aI.e., $X_n(\omega) \rightarrow X(\omega)$ for all $\omega \in \Omega \setminus N$ where $\mathbb{P}(N) = 0$.

However, this might still be too strong.

Definition 2.1.3 (Converge in probability). (X_n) *converges in probability* to X , denoted as $X_n \xrightarrow{p} X$, if for every $\epsilon > 0$, $\mathbb{P}(\|X_n - X\| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.

Remark. $X_n \rightarrow X$ if and only if $\|X_n - X\| \rightarrow 0$. The same also holds for \xrightarrow{p} and $\xrightarrow{\text{a.s.}}$.

A related notion is the following, where we now sum over n .

Definition 2.1.4 (Converge completely). (X_n) *converges completely* to X , denoted as $X_n \xrightarrow{\text{comp}} X$, if for every $\epsilon > 0$, $\sum_{n=1}^{\infty} \mathbb{P}(\|X_n - X\| > \epsilon) < \infty$.

Finally, we have the following.

Definition 2.1.5 (Converge in L^p). (X_n) *converges in L^p* to X for some $p > 0$, denoted as $X_n \xrightarrow{L^p} X$, if $\mathbb{E} [\|X_n - X\|^p] \rightarrow 0$ as $n \rightarrow \infty$.

2.1.1 Connection Between Modes of Convergence

We have the following connections between different modes of convergence.

$$\text{completely} \implies \text{almost-surely} \implies \text{in probability} \longleftarrow \text{in } L^p$$

To show the above, the following characterization for [almost-surely convergence](#) is useful.

Proposition 2.1.1. For a sequence of random vectors (X_n) and a random vector X , we have

$$\begin{aligned} X_n \xrightarrow{\text{a.s.}} X &\Leftrightarrow \mathbb{P}(\|X_k - X\| > \epsilon \text{ for some } k \geq n) \xrightarrow{n \rightarrow \infty} 0 \\ &\Leftrightarrow \mathbb{P}(\|X_n - X\| > \epsilon \text{ for infinitely many } n\text{'s}) = 0 \\ &\Leftrightarrow \mathbb{P}(\limsup_{n \rightarrow \infty} \|X_n - X\| > \epsilon) = 0, \end{aligned}$$

where the above holds for every $\epsilon > 0$.

From [Proposition 2.1.1](#), it's clear that $\xrightarrow{\text{a.s.}}$ implies \xrightarrow{p} since

$$\mathbb{P}(\|X_k - X\| > \epsilon \text{ for some } k \geq n) \geq \mathbb{P}(\|X_n - X\| > \epsilon),$$

hence if the former goes to 0, so does the latter. On the other hand, $\xrightarrow{\text{comp}}$ implies $\xrightarrow{\text{a.s.}}$ follows from the third equivalence. Lastly, the [convergence in \$L^p\$](#) implies the [convergence in probability](#) since

$$\mathbb{P}(\|X_n - X\| > \epsilon) \leq \frac{1}{\epsilon^p} \mathbb{E}[\|X_n - X\|^p]$$

from Markov's inequality. However, the converse is not always true.

Theorem 2.1.1 (Dominated convergence theorem). If $X_n \xrightarrow{p} X$ and $\|X_n - X\| \leq Z$ for all $n \geq 1$ where $\mathbb{E}[\|Z\|^p] < \infty$, then $X_n \xrightarrow{L^p} X$.

Theorem 2.1.2 (Scheffé's theorem). If $X_n \xrightarrow{p} X$ and $\limsup_{n \rightarrow \infty} \mathbb{E}[\|X_n\|^p] \leq \mathbb{E}[\|X\|^p] < \infty$, then $X_n \xrightarrow{L^p} X$.

2.1.2 Consistent Estimator

Let $(X_n) \stackrel{\text{i.i.d.}}{\sim} F$ where F is a distribution function. Say we're interested in some aspect of F , for example, some parameter $\theta = T(F) \in \mathbb{R}^m$. By collecting data X_1, \dots, X_n , we estimate θ by computing an estimator $\hat{\theta}_n$ of θ .¹ There are some properties we might want for $\hat{\theta}_n$.

Definition 2.1.6 (Consistent). $\hat{\theta}_n$ is *consistent* of θ if $\hat{\theta}_n \xrightarrow{p} \theta$ as $n \rightarrow \infty$.

Definition 2.1.7 (Strongly consistent). $\hat{\theta}_n$ is *strongly consistent* of θ if $\hat{\theta}_n \xrightarrow{\text{a.s.}} \theta$ as $n \rightarrow \infty$.

Definition 2.1.8 (Converge in mean squared error). $\hat{\theta}_n$ converges to θ in mean squared error if $\hat{\theta}_n \xrightarrow{L^2} \theta$.

Remark. When $d = 1$, $\mathbb{E}[(\hat{\theta}_n - \theta)^2] = \text{Var}[\hat{\theta}_n] + (\mathbb{E}[\hat{\theta}_n - \theta])^2$. Therefore, $\hat{\theta}_n$ [converges in mean squared error](#) to θ if and only if $\mathbb{E}[\hat{\theta}_n] \rightarrow \theta$ and $\text{Var}[\hat{\theta}_n] \rightarrow 0$.

Let's first see the most well-known estimation problem, the mean estimation.

Example (Mean estimation). Suppose $d = 1$, and let X be non-negative. Say we're interested in $\theta = \mathbb{E}[X]$. It's standard that in this case, we can compute $\mathbb{E}[X]$ by

$$\theta = \mathbb{E}[X] = \int_0^\infty \mathbb{P}(X > t) dt = \int_0^\infty (1 - F(t)) dt.$$

If X has a pmf f , then $\mathbb{E}[X] = \sum_x x f(x) = \sum_x x \Delta F(x)$ where $f(x) = \Delta F(x) \equiv F(x) - F(x^-)$; if

¹ $\hat{\theta}_n$ is a function of X_i 's.

X has a pdf f , then

$$\mathbb{E}[X] = \int_0^\infty xf(x) dx = \int_0^\infty xF(dx).$$

Now, let $\hat{\theta}_n$ to be the sample mean, i.e., $\hat{\theta}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. From the strong law of large number, $\bar{X}_n \xrightarrow{\text{a.s.}} \mathbb{E}[X]$, which implies that $\hat{\theta}_n$ is a **strongly consistent estimator** of θ .

On the other hand, if $\text{Var}[X] < \infty$, then $\bar{X}_n \xrightarrow{L^2} \mathbb{E}[X]$, which further implies $\bar{X}_n \xrightarrow{p} \mathbb{E}[X]$, hence $\hat{\theta}_n$ is **consistent**.^a

^aThe latter is true even without $\text{Var}[X] = \infty$ as we expect.

Proof. We show the last statement. Since $\text{Var}[X] < \infty$, then

$$\frac{\text{Var}[X]}{n} = \text{Var}[\bar{X}_n] = \mathbb{E}[(\bar{X}_n - \mathbb{E}[X])^2] \rightarrow 0$$

as $n \rightarrow \infty$, which implies $\bar{X}_n \xrightarrow{p} \mathbb{E}[X]$. *

Another interesting problem is the supremum estimation.

Example (Supremum estimation). Suppose there is a $\theta \in \mathbb{R}$ where distribution function F such that $F(\theta - \epsilon) < 1 = F(\theta)$ for all $\epsilon > 0$, i.e., $\theta = \sup_{\omega} X(\omega)$ since $\mathbb{P}(X \leq \theta - \epsilon) = F(\theta - \epsilon)$ and $F(\theta) = \mathbb{P}(X \leq \theta)$.^a Then $\hat{\theta}_n = \max_{1 \leq i \leq n} X_i$ is indeed a **strongly consistent** estimator of θ .

^aSuch a distribution exists, for example, $\mathcal{U}(0, \theta)$.

Proof. We see that for any $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon) &= \mathbb{P}(\hat{\theta}_n > \theta + \epsilon) + \mathbb{P}(\hat{\theta}_n < \theta - \epsilon) \\ &= \mathbb{P}\left(\bigcup_{i=1}^n \{X_i > \theta + \epsilon\}\right) + \mathbb{P}\left(\bigcap_{i=1}^n \{X_i < \theta - \epsilon\}\right) \\ &\leq \sum_{i=1}^n \underbrace{\mathbb{P}(X > \theta + \epsilon)}_0 + \prod_{i=1}^n \mathbb{P}(X_i < \theta - \epsilon) = (\mathbb{P}(X_1 < \theta - \epsilon))^n \leq (F(\theta - \epsilon))^n \rightarrow 0 \end{aligned}$$

as $n \rightarrow \infty$ since $F(\theta - \epsilon) < 1$. This shows that $\hat{\theta}_n$ is indeed **consistent**. Moreover, since $\mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon)$ decays exponentially, so this is absolutely summable, hence it's also **strongly consistency**. *

Proving convergence of $\hat{\theta}_n$ is useful, but this might not be enough.

Example. Consider any deterministic sequence (a_n) in \mathbb{R} which converges to 0. Adding a_n to $\hat{\theta}_n$ will not change the convergence of $\hat{\theta}_n$.

The above suggests that we should look at the *distribution* of $\hat{\theta}_n - \theta$ in order to say how does $\hat{\theta}_n \rightarrow \theta$.

Example (Mean estimation for Gaussian). Suppose $X \sim \mathcal{N}(\theta, 1)$. Then $\hat{\theta}_n = \bar{X}_n \sim \mathcal{N}(\theta, 1/n)$, i.e., $\sqrt{n}(\hat{\theta}_n - \theta) \sim \mathcal{N}(0, 1)$. This implies that we can write down a confidence interval (CI) such that $\hat{\theta}_n \pm 1.96/\sqrt{n}$ with 95% CI for $\hat{\theta}_n$.

Doing this for other kind of estimators and F is not that straightforward and will be challenging.

Remark. Let (X_n) and X be d -dimensional random vectors, $h: \mathbb{R}^d \rightarrow \mathbb{R}^m$, and $c \in \mathbb{R}^d$ constant.

(a) If $X_n \rightarrow c$, then $h(X_n) \rightarrow h(c)$ if h is continuous at c .^a This also holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{p} .

(b) If $X_n \rightarrow X$, then $h(X_n) \rightarrow h(X)$ if h is continuous. This also holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{p} .

^aThis is an if and only if condition if this holds for any h .

Let's see some examples.

Example. If $d = 1$, and $X_n \rightarrow \theta \neq 0$. Then $1/X_n \rightarrow 1/\theta$ where

$$h(x) = \begin{cases} \frac{1}{x}, & \text{if } x \neq 0; \\ c, & \text{if } x = 0 \end{cases}$$

for any $c \in \mathbb{R}$. The same holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{p} .

Example. If $X_n \rightarrow X$ and $Y_n \rightarrow Y$, then $(X_n Y_n) \rightarrow (X, Y)$.^a The same holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{p} .

^aThe converse is also true since projections are continuous.

Proof. $\|(X_n, Y_n) - (X, Y)\| \rightarrow 0$ since $\|(X_n, Y_n) - (X, Y)\| \leq \|X_n - X\| + \|Y_n - Y\|$ for all $n \geq 1$.^a The latter two terms goes to 0 (in whatever sense) by assumption. ⊛

^aThis can be seen from $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$.

Lecture 3: Weak Convergence Portmanteau Theorem

2.2 Weak Convergence

25 Jan. 9:30

All convergences we have discussed are in some senses “point-wise” but not “distribution-wise”, and the latter is more powerful. Consider working with a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and the following.

Definition 2.2.1 (Total variation). The *total variation* distance between X and Y in Ω is defined as

$$\text{TV}(X, Y) = \sup_{B \in \mathcal{F}} |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)|$$

Returning to our situation, consider a sequence of random variables (X_n) and a random variable X .

Remark. If X_n has density f_n and X has density f , then $\text{TV}(X_n, X) = \frac{1}{2} \int |f_n - f|$.

Definition 2.2.2 (Converge in total variation). (X_n) *converges in total variation* to X , denoted as $X_n \xrightarrow{\text{TV}} X$, if $\text{TV}(X_n, X) \rightarrow 0$ as $n \rightarrow \infty$.

Remark. If X_n and X have densities f_n and f , $f_n \rightarrow f$ implies $X_n \xrightarrow{\text{TV}} X$ from [Scheffé’s theorem](#).

Note. The above could make sense even if X_n is defined on different $(\Omega_n, \mathcal{F}_n, \mathbb{P}_n)$ for every n .

Let’s see some examples.

Example. Consider $X_n \sim \text{Bin}(n, p_n)$ such that $np_n \rightarrow \lambda \in \mathbb{R}$. As this happens,

$$X_n \sim \text{Bin}(n, p_n) \xrightarrow{\text{TV}} X \sim \text{Pois}(\lambda).$$

Example. Let $X_n \sim f_{\theta_n}$ where $f_{\theta}(x) = f(x)e^{\theta x - \psi(\theta)}$ for some $\theta \in \Theta$. If $\theta_n \rightarrow \theta$, then $X_n \xrightarrow{\text{TV}} X \sim f_{\theta}$. For example, if $X_n \sim \text{Pois}(\theta_n)$ and $\theta_n \rightarrow \theta$, then $X_n \xrightarrow{\text{TV}} X \sim \text{Pois}(\theta)$.

However, [convergence in total variation](#) might be too strong to work with.

Example. Let $X_n \sim \mathcal{U}\{0, 1/n, \dots, (n-1)/n\}$, which should be converging to $X \sim \mathcal{U}(0, 1)$. However, this doesn't happen in total variation distance as we can take B to be \mathbb{Q} .

This suggests that we should look at something weaker.

Definition 2.2.3 (Converge weakly). (X_n) converges weakly to X , denoted as $X_n \xrightarrow{w} X$, if for all bounded continuous $g: \mathbb{R}^d \rightarrow \mathbb{R}$, $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$.

To see how is [weak convergence](#) compared to [convergence in total variation](#), we revisit the above.

Example. Let $X_n \sim \mathcal{U}\{0, 1/n, \dots, (n-1)/n\}$, which should be converging to $X \sim \mathcal{U}(0, 1)$. We have

$$\mathbb{E}[g(X_n)] = \sum_{k=0}^{n-1} g(k/n) \left(\frac{k+1}{n} - \frac{k}{n} \right) \rightarrow \int_0^1 g(x) dx = \mathbb{E}[g(X)]$$

as g is bounded and continuous on $[0, 1]$, hence Riemann integrable.

2.2.1 Portmanteau Theorem

The following is our main tool of proving [weak convergence](#).

Theorem 2.2.1 (Portmanteau theorem). The following are equivalent.

- (a) $X_n \xrightarrow{w} X$.
- (b) $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$ for all bounded Lipschitz $g: \mathbb{R}^d \rightarrow \mathbb{R}$.
- (c) $\mathbb{P}(X \in A) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A)$ for all $A \subseteq \mathbb{R}^d$ open.
- (d) $\mathbb{P}(X \in A) \geq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in A)$ for all $A \subseteq \mathbb{R}^d$ closed.
- (e) $\mathbb{P}(X_n \in A) \rightarrow \mathbb{P}(X \in A)$ for all A such that $\mathbb{P}(X \in \partial A) = 0$.

Before we prove [Portmanteau theorem](#), we should note that all our discussion can be extended to metric spaces from Euclidean spaces. Let's first recall some basic results for metric spaces.

Claim. Given a metric space (S, ρ) , $\rho(\cdot, A)$ is Lipschitz for any $A \subseteq S$, i.e., for any $x, y \in S$,

$$|\rho(x, A) - \rho(y, A)| \leq \rho(x, y).$$

Proof. Since for any $z \in S$, $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$, hence $\rho(x, A) - \rho(y, A) \leq \rho(x, y)$ by taking the infimum over $z \in A$. Interchanging x and y gives another inequality. \circledast

Claim. Given a metric space (S, ρ) , for any $A \subseteq S$, $x \in \overline{A} \Leftrightarrow \rho(x, A) = 0$.

Proof. If $x \in \overline{A}$, there exists (x_n) in A such that $\rho(x_n, x) \rightarrow 0$. Then for any $z \in A$, $\rho(x, z) \leq \rho(x, x_n) + \rho(x_n, z)$, implying

$$\rho(x, A) \leq \rho(x, x_n) + \rho(x_n, A) \rightarrow 0,$$

hence $\rho(x, A) = 0$. On the other hand, suppose $\rho(x, A) = 0$. As $\rho(x, A) = \inf_{y \in A} \rho(x, y)$, there exists (y_n) in A such that $\rho(x, y_n) \rightarrow \rho(x, A) = 0$, i.e., $x \in \overline{A}$. \circledast

The crucial lemma we're going to use to prove [Portmanteau theorem](#) is the following.

Lemma 2.2.1. Given a metric space (S, ρ) and let $A \subseteq S$ be a closed subset. Then there exists bounded Lipschitz $g_k: S \rightarrow \mathbb{R}$, decreasing in k such that $g_k(x) \searrow \mathbb{1}_A(x)$.

Proof. Since A is closed, $A = \overline{A}$ and

$$\mathbb{1}_A(x) = \begin{cases} 1, & \text{if } x \in A \Leftrightarrow \rho(x, A) = 0; \\ 0, & \text{if } x \notin A \Leftrightarrow \rho(x, A) > 0. \end{cases}$$

Now, we let

$$g_k(x) = \begin{cases} 0, & \text{if } \rho(x, A) > \frac{1}{k}; \\ 1 - k\rho(x, A), & \text{otherwise;} \end{cases} = 1 - (k\rho(x, A) \wedge 1).$$

We see that

- if $x \in A$: $\mathbb{1}_A(x) = 1$, and $g_k(x) = 1$ since $\rho(x, A) = 0$;
- if $x \notin A$: $\mathbb{1}_A(x) = 0$, and $\rho(x, A) > 0$ since A closed, and $g_k(x) = 0$ for all large enough k .

Finally, it's clear that $g_k(x)$ takes values in $[0, 1]$, and we now show it's Lipschitz. We have

$$|g_k(x) - g_k(y)| = |(k\rho(x, A) \wedge 1) - (k\rho(y, A) \wedge 1)| \leq k\rho(x, y)$$

for all $x, y \in S$. ■

Then we can prove the [Portmanteau theorem](#).

Proof of Theorem 2.2.1. (a) \Rightarrow (b) is clear. And we start by proving (c) \Leftrightarrow (d).

Claim. (c) \Leftrightarrow (d).

Proof. We first prove that (d) \Rightarrow (c). Since when A is open,

$$\begin{aligned} \mathbb{P}(X \in A) &= 1 - \mathbb{P}(X \in A^c) \leq 1 - \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in A^c) \\ &= 1 - \limsup_{n \rightarrow \infty} (1 - \mathbb{P}(X_n \in A)) = \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A). \end{aligned} \tag{d}$$

(c) \Rightarrow (d) is exactly the same, hence (c) \Leftrightarrow (d). ⊗

Next, we prove (b) \Rightarrow (d), which gives us (a) \Rightarrow (b) \Rightarrow (d) \Leftrightarrow (c).

Claim. (b) \Rightarrow (d).

Proof. From [Lemma 2.2.1](#), there exists bounded Lipschitz $g_k \searrow \mathbb{1}_A$ such that for all closed A ,

$$\mathbb{P}(X_n \in A) = \mathbb{E}[\mathbb{1}_A(X_n)] \leq \mathbb{E}[g_k(X_n)].$$

This is true for every k and n since $g_k \geq \mathbb{1}_A$, and by taking the limit as $n \rightarrow \infty$,

$$\limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in A) \leq \limsup_{n \rightarrow \infty} \mathbb{E}[g_k(X_n)] = \mathbb{E}[g_k(X)]$$

from our assumption (b). Finally, as $k \rightarrow \infty$, it goes to $\mathbb{E}[\mathbb{1}_A(X)] = \mathbb{P}(X \in A)$ as desired. ⊗

The proof will be continued...

Lecture 4: Continuous Mapping Theorem

Before finishing the proof of [Portmanteau theorem](#), we need one additional tool.

30 Jan. 9:30

Lemma 2.2.2. If $\{A_i\}_{i \in I}$ are pairwise disjoint events, then $\{i \in I : \mathbb{P}(A_i) > 0\}$ is countable.^a

^aNote that I can be uncountable.

Proof. Since we can write

$$\{i \in I: \mathbb{P}(A_i) > 0\} = \bigcup_{k=1}^{\infty} \left\{i \in I: \mathbb{P}(A_i) \geq \frac{1}{k}\right\} =: \bigcup_{k=1}^{\infty} I_k,$$

hence it suffices to show $|I_k| < \infty$ for any $k \geq 1$. Indeed, for any k , $|I_k| \leq k$. Suppose not. Then there exists a countable $J_k \subseteq I_k$ such that $|J_k| > k$, implying

$$\mathbb{P}\left(\bigcup_{i \in J_k} A_i\right) = \sum_{i \in J_k} \mathbb{P}(A_i) \geq \frac{|J_k|}{k} > 1,$$

which is a contradiction. ■

We now finish the proof of [Portmanteau theorem](#).

Proof of Theorem 2.2.1 (cont.) We already proved (a) \Rightarrow (b) \Rightarrow (d) \Leftrightarrow (c).

Claim. (c) + (d) \Rightarrow (e).

Proof. We see that for any A , $A^o \subseteq A \subseteq \overline{A}$, and from (c),

$$\begin{aligned} \mathbb{P}(X \in A^o) &\leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A^o) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A) \\ &\leq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in A) \leq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in \overline{A}) \leq \mathbb{P}(X \in \overline{A}) \end{aligned}$$

where the last step follows from (d). Finally, since

$$\mathbb{P}(X \in \overline{A}) - \mathbb{P}(X \in A^o) = \mathbb{P}(\{X \in \overline{A}\} \setminus \{X \in A^o\}) = \mathbb{P}(X \in (\overline{A} \setminus A^o)) = \mathbb{P}(X \in \partial A),$$

which is 0 by our assumption, i.e., inequalities above are all equalities. In particular, since

$$\liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A) \leq \lim_{n \rightarrow \infty} \mathbb{P}(X_n \in A) \leq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in A)$$

$$\text{and } \mathbb{P}(X \in A^o) \leq \mathbb{P}(X \in A) \leq \mathbb{P}(X \in \overline{A}), \quad \mathbb{P}(X \in A) = \lim_{n \rightarrow \infty} \mathbb{P}(X_n \in A). \quad \textcircled{*}$$

Finally, we prove the following.

Claim. (e) \Rightarrow (a).

Proof. For every $g: \mathbb{R}^d \rightarrow \mathbb{R}$ bounded and continuous, we want to show $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$. Suppose $g \geq 0$,^a and let $K \geq g(x)$ for every $x \in \mathbb{R}^d$ (which exists since g is bounded), then

$$\mathbb{E}[g(X_n)] = \int_0^K \mathbb{P}(g(X_n) > t) dt, \quad \mathbb{E}[g(X)] = \int_0^K \mathbb{P}(g(X) > t) dt,$$

so we just need to prove the convergence of the above two integrals. From [bounded convergence theorem](#), it suffices to show that for almost every $t \in [0, K]$,

$$\mathbb{P}(g(X_n) > t) \rightarrow \mathbb{P}(g(X) > t).$$

Observe that $\mathbb{P}(g(X_n) > t) = \mathbb{P}(X_n \in \{g > t\})$ and $\mathbb{P}(g(X) > t) = \mathbb{P}(X \in \{g > t\})$, so from (e) with $A := \{g > t\}$, it suffices to show $\mathbb{P}(X \in \partial\{g > t\}) = 0$ for almost all t . Firstly,

$$\mathbb{P}(X \in \partial\{g > t\}) = \mathbb{P}(X \in \overline{\{g > t\}} \setminus \{g > t\}^o) = \mathbb{P}(X \in \overline{\{g \geq t\}} \setminus \{g > t\}) = \mathbb{P}(g(X) = t).$$

Moreover, consider the events $\{g(X) = t\}_{t \in [0, K]}$, which are pairwise disjoint, hence [Lemma 2.2.2](#) implies $\mathbb{P}(g(X) = t) = 0$ for all but countably many t 's, exactly what we want to show. ⊗

^aOtherwise, we consider $g = g^+ - g^-$ where $g^+ = \max(g, 0)$ and $g^- = \max(-g, 0)$, and everything follows.

This finishes the proof. ■

2.2.2 Continuous Mapping Theorem

A common scenario is that given a nice function h (in terms of continuity), if $X_n \xrightarrow{w} X$, we want to know when will $h(X_n) \xrightarrow{w} h(X)$. To develop the theorem of this, we need some more facts about metric spaces.

As previously seen. Given two metric spaces (S, ρ) , (S', ρ') , $g: S \rightarrow S'$ is continuous if $x_n \xrightarrow{\rho} x$ implies $g(x_n) \xrightarrow{\rho'} g(x)$, or for open $A \subseteq S'$, $g^{-1}(A) \subseteq S$ is open.

Notation. We sometimes write $g^{-1}(A) =: \{g \in A\}$.

It's clear that the following holds.

Note. If $g: S \rightarrow S'$ is continuous and $A \subseteq S'$ is closed, then $\overline{\{g \in A\}} = \{g \in \overline{A}\}$.

However, when g is not continuous and A is not closed, the situation is a bit more complicated. But at least we can first look at the set where g is continuous.

Notation (Continuous set). For any $g: S \rightarrow S'$, we denote the *continuous set* as $C_g := \{x \in S: g \text{ is continuous at } x\}$.

Then we have the following.

Proposition 2.2.1. Given $g: S \rightarrow S'$ between metric spaces and $A \subseteq S'$,

$$C_g \cap \overline{\{g \in A\}} \subseteq \{g \in \overline{A}\}.$$

Proof. Let $x \in C_g \cap \overline{\{g \in A\}}$. Since $x \in \overline{\{g \in A\}}$, there exists $(x_n) \in \{g \in A\}$ such that $x_n \xrightarrow{\rho} x$. Moreover, $x \in C_g$ implies g is continuous at x , hence $g(x_n) \xrightarrow{\rho'} g(x)$, i.e., $g(x) \in \overline{A}$. ■

This allows us to prove the following theorem, which answers our main question in this section.

Theorem 2.2.2 (Continuous mapping theorem). Consider $X_n \xrightarrow{w} X$ and $h: \mathbb{R}^d \rightarrow \mathbb{R}^m$. If $\mathbb{P}(X \in C_h) = 1$, then $h(X_n) \xrightarrow{w} h(X)$.

Proof. Let $A \subseteq \mathbb{R}^m$ be a closed set. Then from [Portmanteau theorem \(d\)](#), we need to show

$$\limsup_{n \rightarrow \infty} \mathbb{P}(h(X_n) \in A) \leq \mathbb{P}(h(X) \in A).$$

Since $\limsup_{n \rightarrow \infty} \mathbb{P}(h(X_n) \in A) = \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in \{h \in A\})$, implying

$$\limsup_{n \rightarrow \infty} \mathbb{P}(h(X_n) \in A) \leq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in \overline{\{h \in A\}}) \leq \mathbb{P}(X \in \overline{\{h \in A\}}),$$

where the last inequality follows again from [Portmanteau theorem \(d\)](#) since $\overline{\{h \in A\}}$ is clearly closed and $X_n \xrightarrow{w} X$. Finally, as $\mathbb{P}(X \in C_h) = 1$,

$$\mathbb{P}(X \in \overline{\{h \in A\}}) = \mathbb{P}(X \in \overline{\{h \in A\}} \cap C_h) \leq \mathbb{P}(X \in \{h \in \overline{A}\})$$

from [Proposition 2.2.1](#), i.e.,

$$\limsup_{n \rightarrow \infty} \mathbb{P}(h(X_n) \in A) \leq \mathbb{P}(X \in \{h \in \overline{A}\}) = \mathbb{P}(X \in \{h \in A\}) = \mathbb{P}(h(X) \in A)$$

since A is closed, hence we're done. ■

Example. Let $d = 1$ and $X_n \xrightarrow{w} X$ where X is continuous. Then $1/X_n \xrightarrow{w} 1/X$ and $X_n^2 \xrightarrow{w} X^2$.

Proof. For the case of $X^2 \xrightarrow{w} X^2$, [continuous mapping theorem](#) clearly applies with $h(x) = x^2$. For the first case, consider

$$h(x) = \begin{cases} \frac{1}{x}, & \text{if } x \neq 0; \\ 0, & \text{if } x = 0. \end{cases}$$

This means $C_h = \mathbb{R} \setminus \{0\}$. Then, we just need to show $\mathbb{P}(X \in C_h) = 1$ and apply [continuous mapping theorem](#). Observe that this is the same as asking $\mathbb{P}(X = 0) = 0$, which is true when X is continuous.^a ⊗

^aEven if X is not continuous, as long as this is true we can conclude the same thing.

Another useful theorem for proving [weak convergence](#) is the following.

Theorem 2.2.3 (Converging together). Let $X_n \xrightarrow{w} X$, and if Y_n on the same probability space as X_n such that $\|X_n - Y_n\| \xrightarrow{p} 0$, i.e., for all $\epsilon > 0$, $\mathbb{P}(\|X_n - Y_n\| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$. Then, $Y_n \xrightarrow{w} X$.

We first see some applications.

Corollary 2.2.1. If $Y_n \xrightarrow{p} X$, then $Y_n \xrightarrow{w} X$. The converse holds as long as $\mathbb{P}(X = c) = 1$ for some constant c .

Proof. By considering $X_n = X$ for all n , [Theorem 2.2.3](#) implies that if $Y_n \xrightarrow{p} X$, $Y_n \xrightarrow{w} X$. Conversely, if $Y_n \xrightarrow{w} c$, from [Portmanteau theorem \(c\)](#), for any fixed $\epsilon > 0$,^a

$$\underbrace{\mathbb{P}(c \in B(c, \epsilon))}_1 \leq \liminf_{n \rightarrow \infty} \mathbb{P}(Y_n \in B(c, \epsilon)),$$

implying $\mathbb{P}(Y_n \in B(c, \epsilon)) \rightarrow 1$, i.e., $\mathbb{P}(\|Y_n - c\| < \epsilon) \rightarrow 1$. ■

^aRecall that $B(c, \epsilon)$ is the open ball centered at c with radius ϵ .

Lecture 5: Convergence in Distribution and Weak Convergence

Now we prove [Theorem 2.2.3](#).

Proof. From [Portmanteau theorem \(b\)](#), we want to prove that $\mathbb{E}[g(Y_n)] \rightarrow \mathbb{E}[g(X)]$ for all bounded and Lipschitz $g: \mathbb{R}^d \rightarrow \mathbb{R}$. Specifically, let $|g(x)| \leq C$ for all $x \in \mathbb{R}^d$ and $|g(x) - g(y)| \leq K\|x - y\|$ for all $x, y \in \mathbb{R}^d$. From triangle inequality,

$$|\mathbb{E}[g(Y_n)] - \mathbb{E}[g(X)]| \leq |\mathbb{E}[g(Y_n)] - \mathbb{E}[g(X_n)]| + |\mathbb{E}[g(X_n)] - \mathbb{E}[g(X)]|.$$

Since $X_n \xrightarrow{w} X$, the second term goes to 0. As for the first term, since Y_n and X_n are in the same probability space, we see that

$$\begin{aligned} |\mathbb{E}[g(Y_n)] - \mathbb{E}[g(X_n)]| &= |\mathbb{E}[g(Y_n) - g(X_n)]| \\ &\leq \mathbb{E}[|g(Y_n) - g(X_n)|] \\ &= \mathbb{E}[|g(Y_n) - g(X_n)| \cdot \mathbb{1}_{\|X_n - Y_n\| > \epsilon}] + \mathbb{E}[|g(Y_n) - g(X_n)| \cdot \mathbb{1}_{\|X_n - Y_n\| \leq \epsilon}] \\ &\leq 2C\mathbb{P}(\|X_n - Y_n\| > \epsilon) + K\epsilon\mathbb{P}(\|X_n - Y_n\| \leq \epsilon) \\ &\leq 2C\mathbb{P}(\|X_n - Y_n\| > \epsilon) + K\epsilon. \end{aligned}$$

As $n \rightarrow \infty$, we finally have

$$\limsup_{n \rightarrow \infty} |\mathbb{E}[g(Y_n)] - \mathbb{E}[g(X)]| \leq K\epsilon$$

for all $\epsilon > 0$, by letting $\epsilon \rightarrow 0$, we're done. ■

We can now apply [Theorem 2.2.3](#) to prove something similar as we have seen before in the case of [convergence in probability](#).

1 Feb. 9:30

As previously seen. $X_n \xrightarrow{p} X$ and $Y_n \xrightarrow{p} Y$ if and only if $(X_n, Y_n) \xrightarrow{p} (X, Y)$.

Now, in the case of **weak convergence**, from **continuous mapping theorem**, we see that if $(X_n, Y_n) \xrightarrow{w} (X, Y)$, then $X_n \xrightarrow{w} X$ and $Y_n \xrightarrow{w} Y$. However, the converse needs not be true.

Example. Consider a random variable X on $(\Omega, \mathcal{F}, \mathbb{P})$, and let $X_n = X$, $Y_n = -X$ for all $n \geq 1$. If $X \sim \mathcal{N}(0, 1)$, we see that $\mathbb{P}(X \in A) = \mathbb{P}(-X \in A)$ for all $A \subseteq \mathbb{R}^d$, implying $X_n \xrightarrow{w} X$ and $Y_n \xrightarrow{w} X$. However, this does not imply $(X_n, Y_n) \xrightarrow{w} (X, X)$ since otherwise, by **continuous mapping theorem**, $X_n + Y_n \xrightarrow{w} X + X = 2X$, which is not true since $X_n + Y_n = 0$.

But in the case of Y is a constant, the converse is actually true, and the result is quite useful.

Theorem 2.2.4 (Slutsky's theorem). If $X_n \xrightarrow{w} X$ in \mathbb{R}^d and $Y_n \xrightarrow{p} c$ in \mathbb{R}^m ,^a then $(X_n, Y_n) \xrightarrow{w} (X, c)$.

^aRecall that from **Corollary 2.2.1**, for a constant c , **weak convergence** is equivalent to **convergence in probability**.

Proof. Firstly, we show that $(X_n, c) \xrightarrow{w} (X, c)$. Indeed, since for every continuous and bounded $g: \mathbb{R}^{d+m} \rightarrow \mathbb{R}$, $\mathbb{E}[g(X_n, c)] \rightarrow \mathbb{E}[g(X, c)]$ follows directly from $X_n \xrightarrow{w} X$ with $g(\cdot, c)$ being continuous and bounded.

Secondly, we show that $\|(X_n, Y_n) - (X_n, c)\| \xrightarrow{p} 0$. This is easy since

$$\|(X_n, Y_n) - (X_n, c)\| \leq \|X_n - X_n\| + \|Y_n - c\| = \|Y_n - c\|,$$

which goes to 0 in probability as we wish. Combining both with **Theorem 2.2.3** gives the result. ■

Revisiting the **counter-example**, we see that now it's not the case when Y is a constant.

Corollary 2.2.2. If $X_n \xrightarrow{w} X$ and $Y_n \xrightarrow{p} c$ in \mathbb{R}^d , $X_n \pm Y_n \xrightarrow{w} X \pm c$, $X_n \cdot Y_n \xrightarrow{w} X \cdot c$. If $d = 1$ and $c \neq 0$, then $X_n/Y_n \xrightarrow{w} X/c$.

Proof. This follows directly from **Slutsky's theorem** and **continuous mapping theorem**. ■

2.2.3 Convergence in Distribution

So far, the notions of convergence we have talked about applies to general probability space, which needs not to be in \mathbb{R}^d in general. However, traditionally, the case in \mathbb{R}^d is considered first.

Intuition. There's a conical ordering available in \mathbb{R}^d to define F_X and F_{X_n} .

This allows us to define the following.

Definition 2.2.4 (Converge in distribution). Let (X_n) and X be random variables in \mathbb{R}^d . Then (X_n) converges in distribution to X , denoted as $X_n \xrightarrow{D} X$, if for all $(t_1, \dots, t_d) \in C_{F_X}$,

$$F_{X_n}(t_1, \dots, t_d) \rightarrow F_X(t_1, \dots, t_d).$$

Note. X_n and X (in \mathbb{R}^d) do not have to be on the same probability space.

Specifically, to see how this relates to what we have seen, recall that

$$F_{X_n}(t_1, \dots, t_d) = \mathbb{P}(X_n^i \leq t_i, \forall 1 \leq i \leq d) = \mathbb{P}(X_n \in (-\infty, t_1] \times \dots \times (-\infty, t_d]),$$

same for F_X . So this reduces to the form we're familiar with, i.e., $\mathbb{P}(X_n \in A)$ for some A .

Remark. $X_n \xrightarrow{TY} X$ implies $X_n \xrightarrow{D} X$.

Proof. Since $X_n \xrightarrow{\text{TV}} X$ means $\mathbb{P}(X_n \in A) \rightarrow \mathbb{P}(X \in A)$ uniformly in A , but $X_n \xrightarrow{D} X$ only requires the above holds for A in the form of $(-\infty, t_1] \times \cdots \times (-\infty, t_d]$, which is weaker. \circledast

There are more classical results that are worth mentioning.

Remark (De Moivre central limit theorem). Let $X_n \sim \text{Bin}(n, p)$, then for every $t \in \mathbb{R}$, as $n \rightarrow \infty$,

$$\mathbb{P}\left(\frac{X_n - np}{\sqrt{np(1-p)}} \leq t\right) \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-u^2/2} du = \Phi(t).$$

Proposition 2.2.2. Let X_n and X be in \mathbb{Z} such that f_n and f are their corresponding pmf's, then

$$f_n \rightarrow f \Leftrightarrow X_n \xrightarrow{\text{TV}} X \Leftrightarrow X_n \xrightarrow{D} X.$$

Proof. The forward implications are clear, so we just need to show $X_n \xrightarrow{D} X$ implies $f_n \rightarrow f$. Since for every $t \in \mathbb{Z}$, since X_n and X are discrete in \mathbb{Z} ,

$$f_n(t) = \mathbb{P}(X_n = t) = \mathbb{P}(X_n \leq t + \epsilon) - \mathbb{P}(X_n \leq t - \epsilon)$$

for some $\epsilon > 0$ small enough. Now, as $t \pm \epsilon$ are in C_X clearly, $X_n \xrightarrow{D} X$ implies

$$\mathbb{P}(X_n \leq t + \epsilon) \rightarrow \mathbb{P}(X \leq t + \epsilon),$$

and the same holds for $t - \epsilon$, hence

$$f_n(t) = \mathbb{P}(X_n = t) = \mathbb{P}(X_n \leq t + \epsilon) - \mathbb{P}(X_n \leq t - \epsilon) \rightarrow \mathbb{P}(X \leq t + \epsilon) - \mathbb{P}(X \leq t - \epsilon) = \mathbb{P}(X = t) = f(t).$$

As this holds for every $t \in \mathbb{Z}$, we're done. \blacksquare

Now, the problem one might have is the following.

Problem. Why not defined for all $t \in \mathbb{R}^d$, rather than $t \in C_{F_X}$?

Answer. Consider for $d = 1$ with $X = c \in \mathbb{R}$, i.e., F_X is the step function at c . To show $X_n \xrightarrow{D} c$, we don't have to show $\mathbb{P}(X_n \leq c) \rightarrow \mathbb{P}(X \leq c) = 1$. Otherwise, if we need to show this for all t , in particular, c , $X_n = c + 1/n$ would not satisfy this. \circledast

If $X_n \xrightarrow{D} X$ and X is continuous, then F_{X_n} converges to F_X not only point-wise, but uniformly.

Remark (Polya's theorem). If F_X is continuous, $X_n \xrightarrow{D} X$ is equivalent as

$$\sup_{t \in \mathbb{R}^d} |F_{X_n}(t) - F_X(t)| \rightarrow 0.$$

Now we have seen various remarks and clarifications about [convergence in distribution](#), the upshot is that, it is actually just a renaming of [weak convergence](#) in \mathbb{R}^d !

Theorem 2.2.5. Given X_n and X in \mathbb{R}^d , then $X_n \xrightarrow{w} X$ if and only if $X_n \xrightarrow{D} X$.

Proof. We prove for the case of $d = 1$, then it's easy to see the same holds for $d \geq 1$. For the forward direction, we want to show that for all $t \in C_{F_X}$, $\mathbb{P}(X_n \leq t) \rightarrow \mathbb{P}(X \leq t)$. Note that $\mathbb{P}(X \leq t) = \mathbb{P}(X \in (-\infty, t])$ and $\mathbb{P}(X_n \leq t) = \mathbb{P}(X_n \in (-\infty, t])$, hence, from [Portmanteau theorem \(e\)](#) with $A = (-\infty, t]$, $X_n \xrightarrow{w} X$ is equivalently as saying $\mathbb{P}(X_n \leq t) \rightarrow \mathbb{P}(X \leq t)$ if

$$\mathbb{P}(X \in \partial(-\infty, t]) = \mathbb{P}(X \in \{t\}) = \mathbb{P}(X = t)$$

is 0. This is indeed the case since $t \in C_{F_X}$, hence we're done.

To show the backward direction, we need the following lemma.

Lemma 2.2.3. $X_n \xrightarrow{D} X$ if and only if for all $x \in \mathbb{R}^d$,

$$F_X(x^-) \leq \liminf_{n \rightarrow \infty} F_{X_n}(x^-) \leq \liminf_{n \rightarrow \infty} F_{X_n}(x) \leq \limsup_{n \rightarrow \infty} F_{X_n}(x) \leq F_X(x).$$

Proof. The backward direction is clear, so we prove the forward direction. When $x \in C_{F_X}$, we're clearly done, so consider $x \notin C_{F_X}$. Firstly, note that $|C_{F_X}^c|$ is countable, so there exists $(x_k) \nearrow x$ and $(y_k) \searrow x$, both in C_{F_X} . Hence, for all $n \geq 1$ and $k \geq 1$,

$$F_{X_n}(x_k) \leq F_{X_n}(x) \leq F_{X_n}(y_k)$$

as F_{X_n} is increasing. We now have for every $k \geq 1$,

$$\begin{aligned} F_X(x_k) &= \lim_{n \rightarrow \infty} F_{X_n}(x_k) && x_k \in C_{F_X} \\ &\leq \liminf_{n \rightarrow \infty} F_{X_n}(x^-) \\ &\leq \liminf_{n \rightarrow \infty} F_{X_n}(x) && F_{X_n} \text{ is increasing} \\ &\leq \limsup_{n \rightarrow \infty} F_{X_n}(x) \\ &\leq \limsup_{n \rightarrow \infty} F_{X_n}(y_k) = F_X(y_k). && y_k \in C_{F_X} \end{aligned}$$

By taking $k \rightarrow \infty$, $F_X(x_k) \rightarrow F_X(x^-)$, while $F_X(y_k) \rightarrow F_X(x)$,^a and we're done.

^aRecall that the distribution function is always right-continuous.

The proof will be *continued*...

Lecture 6: Stochastic Boundedness and Delta Theorem

Before we finish the proof of [Theorem 2.2.5](#), we recall one important characterization of \liminf .

2 Feb. 17:30

As previously seen. Given two real sequence x_n and y_n ,

$$\liminf_{n \rightarrow \infty} (x_n + y_n) \geq \liminf_{n \rightarrow \infty} x_n + \liminf_{n \rightarrow \infty} y_n,$$

where the equality holds when either x_n or y_n converges (not if and only if).

We can then finish the proof of [Theorem 2.2.5](#).

Proof of Theorem 2.2.5 (cont.) Now we can prove the backward direction. From [Portmanteau theorem \(c\)](#), it suffices to show that for every open $A \subseteq \mathbb{R}$, we have

$$\mathbb{P}(X \in A) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A).$$

From the elementary analysis, we see that it suffices to show when $A = (a, b)$ since when $A \subseteq \mathbb{R}$ is open, one can write $A = \bigcup_{k=1}^{\infty} (a_k, b_k)$ where (a_k, b_k) 's disjoint, and have

$$\begin{aligned} \mathbb{P}(X \in A) &= \sum_{k=1}^{\infty} \mathbb{P}(X \in (a_k, b_k)) \\ &\leq \sum_{k=1}^{\infty} \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in (a_k, b_k)) && \text{assume true for each } (a_k, b_k) \\ &\leq \liminf_{n \rightarrow \infty} \sum_{k=1}^{\infty} \mathbb{P}(X_n \in (a_k, b_k)) = \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A), \end{aligned}$$

where the last inequality follows from an induction on $\liminf_{n \rightarrow \infty} (x_n + y_n) \geq \liminf_{n \rightarrow \infty} x_n + \liminf_{n \rightarrow \infty} y_n$. Now, we show that $\mathbb{P}(X \in A) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in A)$ when $A = (a, b)$.

Claim. $\mathbb{P}(X \in (a, b)) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in (a, b))$.

Proof. Observe that $\mathbb{P}(X \in (a, b)) = F_X(b^-) - F_X(a)$, with [Lemma 2.2.3](#), we further have

$$\begin{aligned} \mathbb{P}(X \in (a, b)) &= F_X(b^-) - F_X(a) \\ &\leq \liminf_{n \rightarrow \infty} F_{X_n}(b^-) - \left(\limsup_{n \rightarrow \infty} F_{X_n}(a) \right) \\ &\leq \liminf_{n \rightarrow \infty} F_{X_n}(b^-) + \liminf_{n \rightarrow \infty} (-F_{X_n}(a)) \\ &\leq \liminf_{n \rightarrow \infty} (F_{X_n}(b^-) - F_{X_n}(a)) = \liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in (a, b)), \end{aligned}$$

which proves the claim. \otimes

This proves the case of $d = 1$. \blacksquare

[Theorem 2.2.5](#) means that when talking about random vectors, we can use every result we have proved for the case of [weak convergence](#). Let's see one application, which uses [weak convergence](#)'s result but now prove something about the distribution.

Proposition 2.2.3. If $X_n \xrightarrow{D} X$ and $t_n \rightarrow t \in C_{F_X}$, then $\mathbb{P}(X_n \leq t_n) \rightarrow \mathbb{P}(X \leq t)$.

Proof. We see that from [Corollary 2.2.2](#), $X_n - t_n \xrightarrow{w} X - t$, i.e., $X_n - t_n \xrightarrow{D} X - t$. Hence,

$$\mathbb{P}(X_n \leq t_n) = \mathbb{P}(X_n - t_n \leq 0) = F_{X_n - t_n}(0) \rightarrow F_{X - t}(0) = \mathbb{P}(X - t \leq 0)$$

as long as $0 \in C_{F_{X-t}}$, i.e., $\mathbb{P}(X - t = 0) = \mathbb{P}(X = t) = 0$, which is just $t \in C_{F_X}$ as we assumed. \blacksquare

2.2.4 Stochastic Boundedness

So far we have been talking about the notion of convergence, now we switch the gear a bit and consider boundedness. In this section, let $(X_i)_{i \in I}$ be a family of d -dimensional random vectors defined on probability spaces $(\Omega_i, \mathcal{F}_i, \mathbb{P}_i)$, with the non-empty index set I , which can be either finite or infinite.

Definition 2.2.5 (Bounded in probability). $(X_i)_{i \in I}$ is said to be *bounded in probability* if for every $\epsilon > 0$, there exists an $M > 0$ such that for every $i \in I$,

$$\mathbb{P}(\|X_i\| \geq M) < \epsilon.$$

In other words, for every $\epsilon > 0$, there exists an $M > 0$ such that $\mathbb{P}(\|X_i\| < M) \geq 1 - \epsilon$ for every $i \in I$.

Intuition. For any arbitrary large probability close to 1 we want, one can find an upper-bound M on $\|X_i\|$ uniformly for all $i \in I$.

Note. When $X_i = X$ for every $i \in I$, $(X_i)_{i \in I}$ is trivially [bounded in probability](#).

Proof. Since if not, there exists $\epsilon > 0$, for every $M > 0$, $\mathbb{P}(\|X\| \geq M) \geq \epsilon$. Then as $M \rightarrow \infty$, $\mathbb{P}(\|X\| = \infty) \geq \epsilon$, which is a contradiction since $\|X\| = \infty$. \otimes

Remark. When I is finite, $(X_i)_{i \in I}$ is also trivially [bounded in probability](#). On the other hand, when I is infinite, by considering $X_n = n$ (deterministic), which is not [bounded in probability](#) anymore.

We now provide some sufficient conditions for being [bounded in probability](#).

Proposition 2.2.4. If $(X_i)_{i \in I}$ is bounded in L^p for some $p > 0$, i.e., $\sup_{i \in I} \mathbb{E}[\|X_i\|^p] < \infty$, then $(X_i)_{i \in I}$ is **bounded in probability**.

Proof. Denote $K := \sup_{i \in I} \mathbb{E}[\|X_i\|^p] < \infty$. Since for any $\epsilon > 0$, from Markov's inequality,

$$\mathbb{P}(\|X_i\| > M) \leq \frac{\mathbb{E}[\|X_i\|^p]}{M^p} \leq \frac{K}{M^p} =: \epsilon$$

for $M := \sqrt[p]{K/\epsilon}$. Hence, we're done. \blacksquare

We can generalize some relations between convergence and boundedness from the elementary analysis.

As previously seen. If a deterministic sequence in \mathbb{R} converges, then it's bounded.

In our context, we might expect something like “if $X_n \xrightarrow{p} X$, then (X_n) is **bounded in probability**.” In fact, we have the following “stronger” result where we only require **convergence in distribution**.

Proposition 2.2.5. If $X_n \xrightarrow{D} X$, then (X_n) is **bounded in probability**.

Proof. Fix an $\epsilon > 0$. There is an $M > 0$ such that $\mathbb{P}(\|X\| \geq M) < \epsilon$ since this is a single random vector. To relate this back to X_n , from **Portmanteau theorem (d)**,

$$\epsilon > \mathbb{P}(\|X\| \geq M) = \mathbb{P}(X \in B^c(0, M)) \geq \limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in B^c(0, M)) = \limsup_{n \rightarrow \infty} \mathbb{P}(\|X_n\| \geq M).$$

In other words, $\liminf_{n \rightarrow \infty} \mathbb{P}(\|X_n\| < M) > 1 - \epsilon$, hence there exists an n_0 such that for every $n \geq n_0$, $\mathbb{P}(\|X_n\| < M) \geq 1 - \epsilon$. As for those $n < n_0$, since $\{X_n : n < n_0\}$ is a finite family, we can find $M' > 0$ such that $\mathbb{P}(\|X_n\| < M') > 1 - \epsilon$ for every $n < n_0$. Finally, by considering $M'' := \max(M, M')$, we have $\mathbb{P}(\|X_n\| < M'') > 1 - \epsilon$, i.e., $\mathbb{P}(\|X_n\| \geq M'') < \epsilon$ as desired. \blacksquare

A kind of converse theorem is called **Prokhorov's theorem**, but we won't prove it here right now. We now see another useful characterization that generalizes our intuition in \mathbb{R} . Recall the following.

As previously seen. In \mathbb{R} , if $a_n \rightarrow 0$ and b_n is bounded, $a_n b_n \rightarrow 0$.

The generalization is the following.

Proposition 2.2.6. Let $d = 1$ such that X_n and Y_n are defined on the same probability space. If $X_n \xrightarrow{p} 0$ and Y_n is **bounded in probability**, then $X_n Y_n \xrightarrow{p} 0$.

Proof. Fix an $\epsilon > 0$. We want to show that $\mathbb{P}(|X_n Y_n| > \epsilon) \rightarrow 0$. This is because

$$\begin{aligned} \mathbb{P}(|X_n Y_n| > \epsilon) &= \mathbb{P}(|X_n Y_n| > \epsilon, |Y_n| > M) + \mathbb{P}(|X_n Y_n| > \epsilon, |Y_n| \leq M) \\ &\leq \mathbb{P}(|Y_n| > M) + \mathbb{P}(|X_n Y_n| > \epsilon, |Y_n| \leq M) \leq \mathbb{P}(|Y_n| > M) + \mathbb{P}(|X_n| > \epsilon/M) \end{aligned}$$

for any M . Now, we see that

- since Y_n is **bounded in probability**, there's an $M > 0$ such that $\mathbb{P}(|Y_n| > M) < \epsilon$ for all n ;
- since $X_n \xrightarrow{p} 0$, for the M (depends on the fixed ϵ) above, $\mathbb{P}(|X_n| > \epsilon/M) \rightarrow 0$ as $n \rightarrow \infty$.

We see that the second term always goes to 0, while the first term can always be upper-bounded by ϵ . Hence, by letting $\epsilon \rightarrow 0$, we're done. \blacksquare

We often write the following.

Notation. We write $X_n = o_p(1)$ for $X_n \xrightarrow{p} 0$, and $X_n = O_p(1)$ when (X_n) is **bounded in probability**.

Remark. Proposition 2.2.6 means $o_p(1) \times O_p(1) = o_p(1)$.

Let's see one important application which combines the above. Consider an estimator T_n of θ , and a deterministic sequence b_n which goes to ∞ . In this case, we often have

$$b_n(T_n - \theta) \xrightarrow{D} Y.$$

Example. When $X_n \sim \text{Bin}(n, p)$, then for $b_n = \sqrt{n/p(1-p)} \rightarrow \infty$, $T_n = X_n/n$, and $\theta = p$, we have

$$\frac{X_n - np}{\sqrt{np(1-p)}} = \sqrt{\frac{n}{p(1-p)}} \left(\frac{X_n}{n} - p \right) = b_n(T_n - \theta) \rightarrow Y \sim \mathcal{N}(0, 1).$$

This allows us to compute the rate of convergence and the limiting distribution. But what can we say when we care about $g(T_n)$ for a function g ?

Theorem 2.2.6 (Delta method). Let $\theta \in \mathbb{R}^d$, (T_n) and Y be random vectors in \mathbb{R}^d , and $b_n \rightarrow \infty$ be a positive deterministic sequence. If $b_n(T_n - \theta) \xrightarrow{D} Y$, then $T_n \xrightarrow{P} \theta$. Moreover, if $g: \mathbb{R}^d \rightarrow \mathbb{R}^m$ is differentiable at θ , $b_n(g(T_n) - g(\theta)) \xrightarrow{D} \nabla g(\theta)Y$.

Proof. We first observe that $\|b_n(T_n - \theta)\| \in O_p(1)$ since $b_n(T_n - \theta) \xrightarrow{D} Y$, with [continuous mapping theorem](#) and the fact that $\|\cdot\|$ is continuous, $\|b_n(T_n - \theta)\| \xrightarrow{P} \|Y\|$, so $\|b_n(T_n - \theta)\| \in O_p(1)$ by [Proposition 2.2.5](#). With this, as $b_n \rightarrow \infty$, $T_n \xrightarrow{P} \theta$ since

$$\|T_n - \theta\| = \frac{1}{b_n} \|b_n(T_n - \theta)\| = o(1)O_p(1) \xrightarrow{P} 0$$

as $o(1)O_p(1) = o_p(1)$ from [Proposition 2.2.6](#). For the second claim, since g is differentiable at θ ,

$$\frac{g(x) - g(\theta) - \nabla g(\theta)(x - \theta)}{\|x - \theta\|} \rightarrow 0$$

when $x \rightarrow \theta$. Let $r(x) := g(x) - g(\theta) - \nabla g(\theta)(x - \theta)$ for $x \in \mathbb{R}^d$ be the remainder, and consider

$$h(x) = \begin{cases} 0, & \text{if } x = \theta; \\ \frac{r(x)}{\|x - \theta\|}, & \text{if } x \neq \theta, \end{cases}$$

which is continuous at θ . Rewriting everything, we have

$$r(x) = g(x) - g(\theta) - \nabla g(\theta)(x - \theta) = h(x)\|x - \theta\|$$

for every $x \in \mathbb{R}^d$. Now, let $x = T_n$, multiply both sides by b_n , and take the norm, we see that

$$\|b_n(g(T_n) - g(\theta)) - \nabla g(\theta)b_n(T_n - \theta)\| = \|h(T_n)\| \|b_n(T_n - \theta)\|.$$

We observe the following.

Claim. It suffices to show that the right-hand sides goes to 0 in probability.

Proof. Since it implies that $b_n(g(T_n) - g(\theta))$ has the same weak limit as $\nabla g(\theta)b_n(T_n - \theta)$ from [Theorem 2.2.3](#), i.e., $\nabla g(\theta)Y$ from our assumption with [continuous mapping theorem](#). \otimes

It's enough to show $\|h(T_n)\| = o_p(1)$ since we know that $\|b_n(T_n - \theta)\| = O_p(1)$ and $o_p(1)O_p(1) = o_p(1)$ from [Proposition 2.2.6](#). Indeed, as $T_n \xrightarrow{P} \theta$, $h(T_n) \xrightarrow{P} h(\theta) = 0$ again by [continuous mapping theorem](#) with h being continuous at θ . This further implies $\|h(T_n)\| \xrightarrow{P} 0$ as we desired.^a Combining the above, the result follows. \blacksquare

^aThis involves [continuous mapping theorem](#) and [Corollary 2.2.1](#) since $h(\theta) = 0$, a constant (so does its norm).

Hence, we see that the answer to our original question is rather simple: as $b_n(T_n - \theta) \xrightarrow{D} Y$,

$$b_n(g(T_n) - g(\theta)) \xrightarrow{D} \nabla g(\theta) \cdot Y$$

for any differentiable g at θ .

Lecture 7: Skorohod's Representation Theorem

2.2.5 Skorohod's Representation Theorem

6 Feb. 9:30

So far, we have seen the following.



Now, we show an interesting result that we might not expect.

Theorem 2.2.7 (Skorohod's representation theorem). If $X_n \xrightarrow{D} X$, there exists $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$ on which we can define random vectors (Y_n) and Y such that $Y_n \stackrel{D}{=} X_n$ for all n and $Y \stackrel{D}{=} X$, and $\tilde{\mathbb{P}}(Y_n \rightarrow Y) = 1$.

Intuition. We have [convergence in distribution](#) “implies” [almost surely convergence](#).

We want to prove [Skorohod's representation theorem](#) for $d = 1$. To start, say $X \sim F$ on $(\Omega, \mathcal{F}, \mathbb{P})$. We will consider $F^{-1}(p)$, which exists if there exists a unique $t \in \mathbb{R}$ such that $F(t) = p$, then $F^{-1}(p) = t$. However, this is not really practical since in the discrete case, the preimage might not exist; and even if in the continuous F , when F flats out (at $p = 1$), the preimage is not unique.

Definition 2.2.6 (Quantile). A p^{th} quantile of X is defined as any $t \in \mathbb{R}$ such that

$$\mathbb{P}(X \leq t) \geq p \geq \mathbb{P}(X < t).$$

Now, we can define $F^{-1}(p)$ as the smallest [quantile](#).

Definition 2.2.7 (Quantile function). The *quantile function* of $X \sim F$ is defined as

$$F^{-1}(p) = \inf\{t \in \mathbb{R} : F(t) \geq p\}.$$

We sometimes also call F^{-1} as the *generalized inverse* of F .

Remark. $t \geq F^{-1}(p)$ if and only if $F(t) \geq p$; in other words, $t < F^{-1}(p)$ if and only if $F(t) < p$.

One application of F^{-1} is that given any cdf F , we can construct a corresponding random variable.

Remark (Construction of random variable). Let $U \sim \mathcal{U}(0, 1)$ be a uniform random variable on $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$. Then, $F^{-1}(U) =: Y$ is a random variable with cdf F .

Proof. Since for any $t \in \mathbb{R}$,

$$\tilde{\mathbb{P}}(Y \leq t) = \tilde{\mathbb{P}}(F^{-1}(U) \leq t) = \mathbb{P}(U \leq F(t)) = F(t).$$

⊛

Now we can prove [Skorohod's representation theorem](#).

Proof of Theorem 2.2.7. Consider $\tilde{\Omega} = (0, 1)$, and $\tilde{\mathbb{P}}((a, b)) = b - a$ for all $a < b$. Then, we can define $U(p) = p$ for all $p \in \tilde{\Omega}$, i.e., $U \sim \mathcal{U}(0, 1)$. Define $Y_n = F_{X_n}^{-1}(U)$ and $Y = F_X^{-1}(U)$ from the [quantile functions](#). Denote Φ be the cdf of $\mathcal{N}(0, 1)$, and let $Z = \Phi^{-1}(U)$.

It's clear that $Y_n \stackrel{D}{=} X_n$ and $Y \stackrel{D}{=} X$, so we just need to show $\tilde{\mathbb{P}}(Y_n \rightarrow Y) = 1$.

Claim. It's equivalent to $\tilde{\mathbb{P}}(F_{X_n}(Z) < p) \rightarrow \tilde{\mathbb{P}}(F_X(Z) < p)$ for almost all p 's.

Proof. Observe further that $Y_n(p) = F_{X_n}^{-1}(p)$, $Y(p) = F_X^{-1}(p)$, and $Z(p) = \Phi^{-1}(p)$ for all $p \in (0, 1)$. Since for almost all p 's, $Y_n(p) \rightarrow Y(p)$ if and only if $\Phi(Y_n(p)) \rightarrow \Phi(Y(p))$ as Φ is strictly increasing and continuous, or equivalently,

$$\Phi(Y_n(p)) = \tilde{\mathbb{P}}(Z \leq Y_n(p)) \rightarrow \tilde{\mathbb{P}}(Z \leq Y(p)) = \Phi(Y(p)).$$

As Z is continuous, this is equivalent to $\tilde{\mathbb{P}}(Z < Y_n(p)) \rightarrow \tilde{\mathbb{P}}(Z < Y(p))$, i.e.,

$$\tilde{\mathbb{P}}(Z < F_{X_n}^{-1}(p)) \rightarrow \tilde{\mathbb{P}}(Z < F_X^{-1}(p)),$$

which holds if and only if $\tilde{\mathbb{P}}(F_{X_n}(Z) < p) \rightarrow \tilde{\mathbb{P}}(F_X(Z) < p)$.^a ⊗

^aFollows from [the remark](#). Explicitly, firstly, it's equivalent to $\tilde{\mathbb{P}}(Z \geq F_{X_n}^{-1}(p)) \rightarrow \tilde{\mathbb{P}}(Z \geq F_X^{-1}(p))$, and with $\tilde{\mathbb{P}}(Z \geq F_{X_n}^{-1}(p)) = \tilde{\mathbb{P}}(F_{X_n}(Z) \geq p)$ and $\tilde{\mathbb{P}}(Z \geq F_X^{-1}(p)) = \tilde{\mathbb{P}}(F_X(Z) \geq p)$, the result follows.

Now we show $\tilde{\mathbb{P}}(F_{X_n}(Z) < p) \rightarrow \tilde{\mathbb{P}}(F_X(Z) < p)$ for almost all p 's. Since $X_n \stackrel{D}{\rightarrow} X$ means $F_{X_n}(t) \rightarrow F_X(t)$, from [Lemma 2.2.3](#), it further implies $F_{X_n}(t^-) \rightarrow F_X(t^-)$ for all $t \in C_{F_X}$. Note that $\tilde{\mathbb{P}}(Z \in C_{F_X}) = 1$ since there can be only countably many discontinuities of F_X . Hence,

$$\tilde{\mathbb{P}}(F_{X_n}(Z) \rightarrow F_X(Z)) = 1,$$

i.e., [converges almost surely](#), which implies $F_{X_n}(Z) \stackrel{D}{\rightarrow} F_X(Z)$, i.e., for all $p \in C_{F_X}(Z)$

$$\tilde{\mathbb{P}}(F_{X_n}(Z) \leq p) \rightarrow \tilde{\mathbb{P}}(F_X(Z) \leq p),$$

and also $\tilde{\mathbb{P}}(F_{X_n}(Z) < p) \rightarrow \tilde{\mathbb{P}}(F_X(Z) < p)$ from [Lemma 2.2.3](#). Again, as F_X can have only countably many discontinuities, this holds for almost all p 's, which is what we want to show. ■

We now see some applications of [Skorohod's representation theorem](#), where we can obtain relatively simple proofs for several theorems, such as [Theorem 2.2.5](#).

Remark. If $X_n \stackrel{D}{\rightarrow} X$, from [Skorohod's representation theorem](#), we can obtain $Y_n \xrightarrow{\text{a.s.}} Y$ on $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$ such that $X_n \stackrel{D}{=} Y_n$ and $X \stackrel{D}{=} Y$. Then by the [bounded convergence theorem](#), for any bounded and continuous g ,

$$\mathbb{E}[g(X_n)] = \tilde{\mathbb{E}}[g(Y_n)] \rightarrow \tilde{\mathbb{E}}[g(Y)] = \mathbb{E}[g(X)].$$

Another application is to generalize [Fatou's lemma](#).

Proposition 2.2.7 (Fatou's lemma). Let $X_n \stackrel{D}{\rightarrow} X^a$ and $g: \mathbb{R}^d \rightarrow [0, \infty)$ continuous. Then

$$\mathbb{E}[g(X)] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[g(X_n)].$$

^aCan be on different probability spaces.

Proof. Let $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$, from [Skorohod's representation theorem](#), we can construct $Y_n \stackrel{D}{=} X_n$, $Y \stackrel{D}{=} X$, and $Y_n \xrightarrow{\text{a.s.}} Y$, which implies $g(Y_n) \xrightarrow{\text{a.s.}} g(Y)$. From [Fatou's lemma](#) in $d = 1$, $\tilde{\mathbb{E}}[g(Y)] \leq \liminf_{n \rightarrow \infty} \tilde{\mathbb{E}}[g(Y_n)]$. The result then follows directly from

$$\mathbb{E}[g(X)] = \tilde{\mathbb{E}}[g(Y)] \leq \liminf_{n \rightarrow \infty} \tilde{\mathbb{E}}[g(Y_n)] = \liminf_{n \rightarrow \infty} \mathbb{E}[g(X_n)].$$

The following is well-known from real analysis [dominated convergence theorem](#). ■

Theorem 2.2.8. If $X_n \xrightarrow{\text{a.s.}} X$, $g: \mathbb{R}^d \rightarrow \mathbb{R}$ is continuous and $(g(X_n))$ is uniformly integrable^a if and only if $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$.

^aI.e., $\lim_{t \rightarrow \infty} \sup_{n \geq 1} \mathbb{E}[|g(X_n)| \mathbb{1}_{|g(X_n)| \geq t}] = 0$.

If $X_n \xrightarrow{w} X$, then from the definition, we will have $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$ if g is continuous and bounded. We can indeed relax both continuity and boundedness as follows.

Proposition 2.2.8. If $X_n \xrightarrow{w} X$ and $\mathbb{P}(X \in C_g) = 1$ where $g: \mathbb{R}^d \rightarrow \mathbb{R}$ such that $(g(X_n))$ is uniformly integrable, then $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$.

Proof. From $\mathbb{P}(X \in C_g) = 1$ and $X_n \xrightarrow{w} X$, from [continuous mapping theorem](#), $g(X_n) \xrightarrow{w} g(X)$, hence $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$. ■

Remark. [Proposition 2.2.8](#) can be proved with [Skorohod's representation theorem](#) also.

2.2.6 Proving Distributional Convergence

We often want to prove $X_n \xrightarrow{D} X$, which is not efficient if we start from the [definition](#). To get some intuition for potential proof strategies, consider a deterministic sequence (x_n) in a metric space (S, ρ) .

Theorem 2.2.9. $(x_n) \rightarrow x$ if and only if every subsequence of (x_n) has a subsequence that converges to the same limit x .

Proof. The forward direction is clear. For the backward direction, if not, there exists (x_{n_k}) and $\epsilon > 0$ such that $\rho(x_{n_k}, x) \geq \epsilon$ for every $k \geq 1$. But if there exists a subsubsequence $(x_{n_{k_\ell}})$ that converges to x , this is clearly a contradiction. ■

In the same vein, with the same argument, we have the following.

Theorem 2.2.10. $X_n \xrightarrow{w} X$ if and only if every subsequence of (X_n) has a subsequence that [converges weakly](#), and all [weakly convergent](#) subsequences have the same limit X .

Proof. Mimicking the proof as in [Theorem 2.2.9](#). ■

Lecture 8: Characteristic Functions

We see other similar theorems apart from [Theorem 2.2.10](#).

8 Feb. 9:30

Theorem 2.2.11. If $X_n \xrightarrow{w} X$ and $X_n \xrightarrow{w} Y$, then $X \stackrel{D}{=} Y$. More generally, if $X_n \xrightarrow{w} X$ and $Y_n \xrightarrow{w} Y$, with $X_n \stackrel{D}{=} Y_n$ for all $n \geq 1$, $X \stackrel{D}{=} Y$.

Proof. We have for every $n \geq 1$, $\mathbb{E}[g(X_n)] = \mathbb{E}[g(Y_n)]$ for all $g: \mathbb{R}^d \rightarrow \mathbb{R}$. If g is bounded and continuous, $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$ and $\mathbb{E}[g(Y_n)] \rightarrow \mathbb{E}[g(Y)]$. To show that $X \stackrel{D}{=} Y$, we want to show $F_X = F_Y$, or $\mathbb{P}(X \in B) = \mathbb{P}(Y \in B)$ for all $B \in \mathcal{F} = \mathcal{B}(\mathbb{R}^d)$. In fact, it's enough to show this for closed B . With [Lemma 2.2.1](#), there exists $(g_k) \searrow \mathbb{1}_B$ for closed B and bounded, Lipschitz g_k , i.e.,

$$\begin{aligned} \mathbb{E}[\mathbb{1}_B(X)] &= \lim_{k \rightarrow \infty} \mathbb{E}[g_k(X)] = \lim_{k \rightarrow \infty} \lim_{n \rightarrow \infty} \mathbb{E}[g_k(X_n)] \\ &= \lim_{k \rightarrow \infty} \lim_{n \rightarrow \infty} \mathbb{E}[g_k(Y_n)] = \lim_{k \rightarrow \infty} \mathbb{E}[g_k(Y)] = \mathbb{E}[\mathbb{1}_B(Y)], \end{aligned}$$

where the third equality follows from the fact that $X_n \stackrel{D}{=} Y_n$. ■

One question is that, if we don't have things like [weak convergent](#) but just some moment information (i.e., when $g(x) = x^k$ when computing $\mathbb{E}[g(X)]$), can we conclude the same thing?

Problem (Method of Moments). If $\mathbb{E}[X^k] = \mathbb{E}[Y^k] < \infty$ for all $k \geq 1$, does $X \stackrel{D}{=} Y$?

Answer. Not in general. We will discuss this more in the assignment. \circledast

2.3 Characteristic Function

To answer the question left above, we will see that it actually suffices to show only for $g(x) = \cos(t \cdot x)$ or $\sin(t \cdot x)$ for $t, x \in \mathbb{R}^d$. This leads to the so-called **characteristic functions**.

Definition 2.3.1 (Characteristic function). The *characteristic function* of a d -dimensional random vector X is defined as $\phi_X: \mathbb{R}^d \rightarrow \mathbb{C}$ where $t \in \mathbb{R}^d$ such that

$$\phi_X(t) = \mathbb{E}[\cos(t \cdot X)] + i\mathbb{E}[\sin(t \cdot X)] = \mathbb{E}[e^{i(t \cdot X)}].$$

Notation. We will now drop the inner product, i.e., write $t \cdot X =: tX$.

If we write ϕ_X explicitly, we have

$$\phi_X(t) = \mathbb{E}[e^{itX}] = \int e^{itx} f_X(x) dx = \int e^{itx} F_X(dx).$$

Remark. **Characteristic functions** are bounded.

Proof. Since

$$|\phi_X(t)| = \sqrt{(\mathbb{E}[\cos(tX)])^2 + (\mathbb{E}[\sin(tX)])^2} \leq \sqrt{\mathbb{E}[\cos^2(tX)] + \mathbb{E}[\sin^2(tX)]} = 1.$$

\circledast

This implies that ϕ_X is meaningful for any random vector X , unlike the moment generating function.

Remark. If X and Y are independent, $\phi_{X+Y}(t) = \phi_X(t) \cdot \phi_Y(t)$.

We make one more remark for future reference.

Remark. If X, Y are discrete, $f_{X+Y}(x) = \sum_y f_Y(x-y)f_X(y)$. More generally, if X, Y have pdfs,

$$f_{X+Y}(x) = \int f_Y(x-y)f_X(y) dy = \int f_Y(x-y)F_X(dy).$$

Furthermore, even if X doesn't have pdf, as long as Y does, the above still holds.

2.3.1 Uniqueness Theorem

Now we can prove the following uniqueness theorem, which states that indeed, it suffices to check only $\sin(tx)$ and $\cos(tx)$ when proving **weak convergence**.

Theorem 2.3.1 (Uniqueness). If $\phi_X(t) = \phi_Y(t)$ for all $t \in \mathbb{R}^d$, then $X \stackrel{D}{=} Y$. The converse is trivial.

Proof. Consider $d = 1$. Observe that if we can write F_X in terms of only ϕ_X , then $\phi_X = \phi_Y$ implies $F_X = F_Y$. To do this, consider the following.

Claim. For $Z, Z' \sim \mathcal{N}(0, 1)$ (independent of X and Y), if one can write $F_{X+\sigma Z}$ for all $\sigma > 0$ in terms of only ϕ_X , $\phi_X = \phi_Y$ implies $X \stackrel{D}{=} Y$.

Proof. Fix some $\sigma > 0$. In this case, if we can write $F_{X+\sigma Z}$ in terms of only ϕ_X , $\phi_X = \phi_Y$ implies $F_{X+\sigma Z} = F_{Y+\sigma Z'}$. This implies $X + \sigma Z \stackrel{D}{=} Y + \sigma Z'$. Now, for $\sigma = 1/k$, $k \in \mathbb{N}$,

$$X + \frac{1}{k}Z \stackrel{D}{=} Y + \frac{1}{k}Z'.$$

With [Corollary 2.2.2](#), since $Z/k \xrightarrow{P} 0$ (and also $Z'/k \xrightarrow{P} 0$), we have $X + Z/k \xrightarrow{D} X$ and $Y + Z'/k \xrightarrow{D} Y$, which implies $X \stackrel{D}{=} Y$ from [Theorem 2.2.11](#). \otimes

Hence, our goal now is to write $F_{X+\sigma Z}$ in terms of ϕ_X . Firstly, for all $t \in \mathbb{R}$,

$$\phi_Z(t) = \int e^{itz} F_Z(dz) = \int e^{itz} f_Z(z) dz = \int e^{itz} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = e^{-t^2/2}. \quad (2.1)$$

Now, consider $f_{X+\sigma Z}(x)$ instead, which exists since Z has a pdf from the [remark](#). We see that

$$\begin{aligned} f_{X+\sigma Z}(x) &= \int f_{\sigma Z}(x-y) F_X(dy) \\ &= \int \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-y)^2/2\sigma^2} F_X(dy), \end{aligned}$$

by replacing $e^{-(x-y)^2/2\sigma^2}$ from [Equation 2.1](#) with $t = (x-y)/\sigma$,

$$\begin{aligned} &= \int \frac{1}{\sigma\sqrt{2\pi}} \int e^{i\frac{y-x}{\sigma}z} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz F_X(dy). \\ &= \frac{1}{2\pi} \iint e^{i(y-x)u} e^{-\sigma^2 u^2/2} du F_X(dy), \quad z/\sigma =: u \\ &= \frac{1}{2\pi} \int e^{-ixu - \sigma^2 u^2/2} \underbrace{\int e^{iyu} F_X(dy)}_{\phi_X(u)} du, \end{aligned}$$

where we interchange the order of integrals with [Fubini's theorem](#) (justified by [Tonelli's theorem](#)) when integrands are absolute integrable. This implies that $F_{X+\sigma Z}(dx)$ can be written in terms of ϕ_X where with no other dependencies, hence we're done. \blacksquare

Note. Now showing $X \stackrel{D}{=} Y$ reduces to calculus.

2.3.2 Continuity Theorem

One immediate consequence of the [uniqueness theorem](#) is that it's enough to have the [characteristic functions](#) converging to some function (not necessarily a [characteristic functions](#) of some X) for us to conclude that the subsequences of (X_n) have the same weak limit. To do this, we need to prove [Prokhorov's theorem](#).

Theorem 2.3.2 (Prokhorov's theorem). If $(X_n) = O_p(1)$, then there exists a [weakly convergent](#) subsequence of (X_n) .

Proof. Based on [Helly's selection theorem](#), $F_{X_n}(t) \rightarrow F(t)$ for all $t \in C_F$, there exists an increasing F , right continuous, $F(+\infty) \leq 1$ and $F(-\infty) \geq 0$ (called the *defective cdf*). Consider $d = 1$, we show that this F is indeed a cdf when $X_n = O_p(1)$.

Fix $\epsilon > 0$, then there exists $M_\epsilon > 0$ in C_F such that

$$F_{X_n}(M_\epsilon) = \mathbb{P}(X_n \leq M_\epsilon) \geq \mathbb{P}(|X_n| \leq M_\epsilon) \geq 1 - \epsilon$$

for all $n \geq 1$. Since $M_\epsilon \in C_F$, $F_{X_n}(M_\epsilon) \rightarrow F(M_\epsilon)$. We then see that for all $\epsilon > 0$, there exists $M_\epsilon > 0$ such that $F(+\infty) \geq F(M_\epsilon) \geq 1 - \epsilon$. As $\epsilon \rightarrow 0$, $F(+\infty) = 1$. Similarly, $F(-\infty) = 0$. \blacksquare

We now state the theorem.

Theorem 2.3.3 (Lévy-Cramer continuity theorem). If $\phi_{X_n}(t) \rightarrow \phi(t)$ for all $t \in \mathbb{R}^d$, then all **weakly convergent** subsequences of (X_n) have the same weak limit. Furthermore, if also ϕ is continuous at 0, then there exists X such that $\phi = \phi_X$ and $X_n \xrightarrow{D} X$.

Proof. Let's start with the first claim. Suppose $Y_n \xrightarrow{w} Y$ and $Z_n \xrightarrow{w} Z$ are two subsequences of X_n such that $Y \neq Z$. But since $\phi_{Y_n}(t) \rightarrow \phi_Y(t)$ and $\phi_{Z_n}(t) \rightarrow \phi_Z(t)$, with the fact that $(\phi_{Y_n}(t))$ and $(\phi_{Z_n}(t))$ are subsequences of $(\phi_{X_n}(t))$ for every t , as $\phi_{X_n}(t) \rightarrow \phi(t)$, both subsequences need to converge to the same limit, i.e.,

$$\phi_Y(t) = \phi(t) = \phi_Z(t)$$

for all $t \in \mathbb{R}^d$. From the **uniqueness theorem**, $Y \stackrel{D}{=} Z$.

For the second part, we just need to prove the following.

Claim. It's enough to show that if ϕ is continuous at 0, $(X_n) = O_p(1)$.

Proof. Since if we can show $(X_n) = O_p(1)$ from our assumption, **Prokhorov's theorem** implies there exists a **weakly convergent** subsequence of (X_n) . With the first claim, we can find the weak limit X . ⊗

The proof will be **continued**...

Lecture 9: Proof of Lévy-Cramer Continuity Theorem

We now finish the proof of **Lévy-Cramer continuity theorem**.

13 Feb. 9:30

Proof of Theorem 2.3.3 (cont.) Fix $\epsilon > 0$. Then there exists $\delta > 0$ such that for all $|t| < \delta$,

$$|\phi(t) - \phi(0)| = |\phi(t) - 1| < \frac{\epsilon}{4}$$

since for any $n \geq 1$, $\phi_{X_n}(0) = 1$, so is $\phi(0)$. Hence, we have

$$\frac{\epsilon}{2} = \frac{1}{\delta} \int_{-\delta}^{\delta} \frac{\epsilon}{4} dt > \frac{1}{\delta} \int_{-\delta}^{\delta} |\phi(t) - 1| dt.$$

We claim that we can find an $n_0 \in \mathbb{N}$ such that for every $n \geq n_0$, $\mathbb{P}(|X_n| \geq 2/\delta) < \epsilon$.^a To bound $|X_n|$ with ϕ_{X_n} , firstly, for all x , $|\sin x| \leq |x|$. This bound is good only when x is close to 0. If it's not the case, then we can use $|\sin x/x| \leq 1/|x| \leq 1/2$ if $|x| \geq 2$. Hence, in general, for $x \neq 0$,

$$\frac{\sin x}{x} \leq \left| \frac{\sin x}{x} \right| \leq \frac{1}{2} \cdot \mathbb{1}_{|x| \geq 2} + 1 \cdot \mathbb{1}_{|x| < 2} = 1 - \frac{1}{2} \mathbb{1}_{|x| \geq 2} \Rightarrow \mathbb{1}_{|x| \geq 2} \leq 2 \left(1 - \frac{\sin x}{x} \right).$$

as $\mathbb{1}_{|x| < 2} = 1 - \mathbb{1}_{|x| \geq 2}$. Plug in δx , for any $x \neq 0$, we have

$$\mathbb{1}_{|\delta x| \geq 2} \leq 2 \left(1 - \frac{\sin(\delta x)}{\delta x} \right) = \frac{1}{\delta} \left(2\delta - 2 \frac{\sin(\delta x)}{x} \right) = \frac{1}{\delta} \int_{-\delta}^{\delta} 1 - \cos(tx) dt.$$

Indeed, the above is true for all $x \in \mathbb{R}$ by manually checking. Finally, by replacing x by X_n and take the expectation on the both sides,

$$\mathbb{P}(|\delta X_n| \geq 2) \leq \frac{1}{\delta} \int_{-\delta}^{\delta} 1 - \mathbb{E}[\cos(tX_n)] dt = \frac{1}{\delta} \int_{-\delta}^{\delta} \operatorname{Re}(1 - \phi_{X_n}(t)) dt \leq \frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi_{X_n}(t)| dt,$$

where we pass the expectation (i.e., limit) inside the integral from **Fubini's theorem** since $\cos(tX_n)$ is bounded. It remains to show that there is some $\delta > 0$ such that the right-hand side is less than ϵ for all $n \geq n_0$. As $\phi_{X_n}(t) \rightarrow \phi(t)$ for all t , we have $|1 - \phi_{X_n}(t)| \rightarrow |1 - \phi(t)|$ point-wise, hence by

the **bounded convergence theorem**,

$$\frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi_{X_n}(t)| dt \rightarrow \frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi(t)| dt < \frac{\epsilon}{2}$$

from our assumption. Putting everything together, there is an $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$,

$$\mathbb{P}(|\delta X_n| \geq 2) = \mathbb{P}(|X_n| \geq 2/\delta) \leq \frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi_{X_n}(t)| dt < \frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi(t)| dt + \frac{\epsilon}{2} < \epsilon,$$

where the second-last inequality follows from the point-wise convergence of $\frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi_{X_n}(t)| dt$ to $\frac{1}{\delta} \int_{-\delta}^{\delta} |1 - \phi(t)| dt$ being $\epsilon/2$ -close for n large enough, i.e., when $n \geq n_0$ for some n_0 . ■

^aIf this is the case, then we can handle the $n < n_0$ case easily as usual by taking the maximum over all $n < n_0$.

2.3.3 Inversion Theorem

On the other hand, another way to prove **Lévy-Cramer continuity theorem** is to directly calculate the pdf of X , given ϕ_X . It follows the same vein of the proof of **uniqueness theorem**.

Intuition. In the proof of **uniqueness theorem**, we only obtain a pdf for $X + \sigma Z$. Imposing constraints on ϕ_X and calculate $\mathbb{E}[g(X)]$ in terms of ϕ_X will tell us which condition should we add.

Theorem 2.3.4 (Feller's inversion formula). Let X be a d -dimensional random vector with the **characteristic function** ϕ_X .

- (a) If g has a bounded support and $\mathbb{P}(X \in C_g) = 1$, then

$$\mathbb{E}[g(X)] = \lim_{\sigma \searrow 0} \frac{1}{2\pi} \iint g(x) e^{-iux - \sigma^2 u^2/2} du dx.$$

- (b) For any $a, b \in C_{F_X}$,

$$F_X(b) - F_X(a) = \lim_{\sigma \searrow 0} \frac{1}{2\pi} \int_a^b \int e^{-iux - \sigma^2 u^2/2} \phi_X(u) du dx.$$

- (c) If further, ϕ_X is absolute integrable, then X has a pdf

$$f_X(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{-iux} \phi_X(u) du.$$

Proof. The proof is based on **uniqueness theorem**.

- (a) In the **uniqueness theorem**, $\sigma \searrow 0$ such that $X + \sigma Z \xrightarrow{D} X$, which implies $g(X + \sigma Z) \xrightarrow{D} g(X)$ when $\mathbb{P}(X \in C_g) = 1$. Since now g is also bounded, by the **bounded convergence theorem**,

$$\mathbb{E}[g(X)] = \lim_{\sigma \searrow 0} \mathbb{E}[g(X + \sigma Z)].$$

We now calculate $\mathbb{E}[g(X + \sigma Z)]$. Since $g: \mathbb{R} \rightarrow \mathbb{R}$ has bounded support, the same calculation from the proof of **uniqueness theorem** gives

$$\mathbb{E}[g(X + \sigma Z)] = \lim_{\sigma \searrow 0} \frac{1}{2\pi} \int g(x) \int e^{-ixu - \sigma^2 u^2/2} \phi_X(u) du dx.$$

It remains to change the order of integration, which is justified by **Tonelli's theorem** as $\mathbb{E}[|g(X + \sigma Z)|] < \infty$ for all $\sigma > 0$, hence we obtain the result for the first part.

- (b) Given $a, b \in C_{F_X}$, consider $g(x) = \mathbb{1}_{(a,b)}(x)$, which implies $\mathbb{P}(X \in C_g) = 1$ (and trivially g has a bounded support), hence the result above applies.
- (c) Finally, if ϕ_X is absolute integrable, our goal now is to pass the limit $\sigma \searrow 0$ inside the integral for $F_X(b) - F_X(a)$ given $a, b \in C_{F_X}$, i.e., to get

$$F_X(b) - F_X(a) = \frac{1}{2\pi} \int_a^b \int \lim_{\sigma \searrow 0} e^{-iux - \sigma^2 u^2 / 2} \phi_X(u) \, du \, dx = \frac{1}{2\pi} \int_a^b \int e^{-iux} \phi_X(u) \, du \, dx,$$

which will imply the result since a cdf is characterized by its values in C_{F_X} , i.e., if the above equality is true, it will be valid for all $a, b \in \mathbb{R}$. To do so, **dominated convergence theorem** states that

$$\int_a^b \int \sup_{\sigma > 0} |e^{-iux - \sigma^2 u^2 / 2} \phi_X(u)| \, du \, dx < \infty$$

is the right condition. We see that the left-hand side is less than

$$\int_a^b \int_{\mathbb{R}} |\phi_X(u)| \sup_{\sigma > 0} |e^{-\sigma^2 u^2 / 2}| \, du \, dx \leq \int_a^b \int_{\mathbb{R}} |\phi_X(u)| \, du \, dx$$

which is finite since $\int |\phi_X(u)| \, du < \infty$. ■

Corollary 2.3.1. Given X_n and X such that ϕ_X and ϕ_{X_n} are integrable. If $\phi_{X_n} \xrightarrow{L^1} \phi_X$,^a then $X_n \xrightarrow{\text{TV}} X$.

^aI.e., $\int_{\mathbb{R}} |\phi_{X_n}(t) - \phi_X(t)| \, dt \rightarrow 0$.

Proof. It suffices to prove that $|f_{X_n}(x) - f_X(x)| \rightarrow 0$, where these pdfs exist due to **Feller's inversion formula (c)**. We see that

$$|f_{X_n}(x) - f(x)| \leq \frac{1}{2\pi} \int_{\mathbb{R}} |e^{-iux}| \cdot |\phi_{X_n}(u) - \phi_X(u)| \, du \leq \frac{1}{2\pi} \int_{\mathbb{R}} |\phi_{X_n}(u) - \phi_X(u)| \, du$$

with the assumption the right-hand side goes to 0. ■

Finally, we see the following characterizations of ϕ_X . The first one is that it's uniformly continuous.

Proposition 2.3.1. For any random vector X , ϕ_X is uniformly continuous, i.e.,

$$\lim_{h \rightarrow 0} \sup_t |\phi_X(t+h) - \phi_X(t)| = 0.$$

Proof. We see that for any h ,

$$|\phi_X(t+h) - \phi_X(t)| = |\mathbb{E}[e^{i(t+h)X}] - \mathbb{E}[e^{itX}]| \leq \mathbb{E}[|e^{itX}| |e^{ihX} - 1|] \leq \mathbb{E}[|e^{ihX} - 1|],$$

which goes to 0 as $h \rightarrow 0$ since $|e^{ihX} - 1| \leq 2$ with **bounded convergence theorem**. ■

The next theorem gives us a way to calculate the derivatives of ϕ_X .

Theorem 2.3.5. If $X \in L^p$ for any $p \in \mathbb{N}$, then the p^{th} derivative of $\phi_X(t)$ is given by

$$\phi_X^{(p)}(t) = \mathbb{E}[(iX)^p e^{itX}]$$

for every t . In particular, $\phi_X^{(p)}(0) = i^p \mathbb{E}[X^p]$.

Proof. We prove the case $p = 1$. We see that it's enough to prove

$$\lim_{h \rightarrow 0} \left| \frac{\phi_X(t+h) - \phi_X(t)}{h} - \mathbb{E}[iX e^{itX}] \right| = 0$$

Writing the ϕ_X explicitly, by Jensen's inequality, for any $h \neq 0$, the left-hand side is

$$\begin{aligned} \left| \frac{\mathbb{E}[e^{i(t+h)X}] - \mathbb{E}[e^{itX}] - \mathbb{E}[ihX e^{itX}]}{h} \right| &\leq \frac{\mathbb{E}[|e^{i(t+h)X} - e^{itX} - ihX e^{itX}|]}{|h|} \\ &= \frac{\mathbb{E}[|e^{itX}| |e^{ihX} - 1 - ihX|]}{|h|} \leq \frac{\mathbb{E}[|e^{ihX} - 1 - ihX|]}{|h|} \end{aligned}$$

Let $G(h) = e^{ihX}$, then $G'(h) = iX e^{ihX}$, and the right-hand side is equal to

$$\frac{\mathbb{E}[|G(h) - G(0) - G'(0)h|]}{|h|}.$$

Since G is differentiable, $G(h) - G(0) = \int_0^h G'(y) dy$, hence

$$G(h) - G(0) - G'(0)h = \int_0^h G'(y) - G'(0) dy = h \int_0^1 G'(uh) - G'(0) du = h \int_0^1 iX e^{iuhX} - iX du$$

where we let $y = uh$. Plugging in, we have

$$\begin{aligned} \mathbb{E} \left[\frac{|e^{ihX} - 1 - ihX|}{|h|} \right] &\leq \mathbb{E} \left[\int_0^1 |G'(uh) - G'(0)| du \right] \\ &= \mathbb{E} \left[\int_0^1 |iX e^{iuhX} - iX| du \right] \leq \mathbb{E} \left[|X| \int_0^1 |e^{iuhX} - 1| du \right]. \end{aligned}$$

Finally, taking the limit as $h \rightarrow 0$, with the fact that $\mathbb{E}[|X|] < \infty$ and $\int_0^1 |e^{ihuX} - 1| du \leq 2$, we see that $|X| \int_0^1 |e^{ihuX} - 1| du \leq 2|X|$, and the latter is integrable since $\mathbb{E}[|X|] < \infty$, hence **dominated convergence theorem** applies, i.e., we can pass the limit into the expectation,

$$\lim_{h \rightarrow 0} \mathbb{E} \left[|X| \int_0^1 |e^{ihuX} - 1| du \right] = \mathbb{E} \left[|X| \lim_{h \rightarrow 0} \int_0^1 |e^{ihuX} - 1| du \right] = 0$$

since $\lim_{h \rightarrow 0} \int_0^1 |e^{ihuX} - 1| du = 0$, again from the **bounded convergence theorem**. ■

Remark. Theorem 2.3.5 implies $\sup_t |\phi_X^{(p)}(t)| \leq \mathbb{E}[|X|^p]$.

Lecture 10: Law of Large Number and Central Limit Theorem

2.4 Fundamental Theorems of Probability

15 Feb. 9:30

With the tools we developed, we now prove two of the fundamental theorems of probability, i.e., the **weak law of large number**, and also the **central limit theorem**. We start from the first one.

Theorem 2.4.1 (Weak law of large number). Let (X_n) be i.i.d. random vectors, and X be a random vector with $\mathbb{E}[|X|] < \infty$. Then $\bar{X}_n \xrightarrow{P} \mathbb{E}[X]$.

Proof. Since $c := \mathbb{E}[X]$ is a constant, it suffices to show that $\phi_{\bar{X}_n}(t) \rightarrow \phi_c(t) = e^{itc}$ for all t from

Corollary 2.2.1. Firstly, let $\bar{X}_n = S_n/n$, we have

$$\phi_{\bar{X}_n}(t) = \mathbb{E}[e^{itS_n/n}] = \phi_{S_n}(t/n) = \prod_{i=1}^n \phi_{X_i}(t/n) = (\phi(t/n))^n$$

where we let $\phi_{X_i} =: \phi$ since (X_n) are i.i.d. From the fundamental theorem of calculus, with the fact that the first moment of X exists, ϕ is differentiable such that

$$(\phi(t/n))^n = \left(1 + \frac{t}{n} \int_0^1 \phi'(ut/n) du\right)^n.$$

Since $(1+a_n)^n \rightarrow e^c$ if $na_n \rightarrow c$, it remains to show $\int_0^1 \phi'(ut/n) du \rightarrow ic$. First, if $\phi'(t)$ is continuous at 0, as $n \rightarrow \infty$

$$\phi'(ut/n) \rightarrow \phi'(0) = i\mathbb{E}[X] = ic.$$

With the fact that $\sup_t |\phi'(t)| \leq \mathbb{E}[|X|]$, the **bounded convergence theorem** implies

$$\int_0^1 \phi'(ut/n) du \rightarrow \int_0^1 ic du = ic$$

since we can now pass the limit inside the integral. ■

Remark. Actually, we don't need to assume finite first moment since assuming ϕ is differentiable at 0 such that $\phi'(0) = ic$ is enough.

In terms of the distributional result, we need higher-order moments. In particular, if the second moment exists, then we can generalize we have done as in the proof of [Theorem 2.3.5](#).

As previously seen. If g is differentiable, then

$$g(x) = g(0) + g'(0)x + x \int_0^1 g'(ux) - g'(0) du.$$

Note. If g is twice-differentiable,

$$\begin{aligned} g(x) &= g(0) + g'(0)x + x \int_0^1 \int_0^{ux} g''(y) dy du \\ &= g(0) + g'(0)x + x \int_0^1 \int_0^1 g''(xuv) ux dv du & y = xuv, dy = ux dv \\ &= g(0) + g'(0)x + x^2 \int_0^1 \int_0^1 g''(xuv) u dv du. \end{aligned}$$

We now state the theorem.

Theorem 2.4.2 (Central limit theorem). Let (X_n) be i.i.d. random variables (i.e., $d = 1$) with $\mathbb{E}[X_i] =: \mu$, $\text{Var}[X_i] =: \sigma^2 < \infty$ for all $1 \leq i \leq n$. Then

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} \mathcal{N}(0, 1).$$

Proof. Without loss of generality, let $\mu = 0$, $\sigma = 1$. Since $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{S_n - n\mu}{\sigma\sqrt{n}}$, it's enough to show that $\phi_{S_n/\sqrt{n}}(t) \rightarrow e^{-t^2/2}$ for any $t \in \mathbb{R}$ from [Lévy-Cramer continuity theorem](#) and [Equation 2.1](#). Firstly,

$$\phi_{S_n/\sqrt{n}}(t) = \mathbb{E}[e^{itS_n/\sqrt{n}}] = \phi_{S_n}(t/\sqrt{n}) = (\phi(t/\sqrt{n}))^n$$

where we let $\phi_{X_n} =: \phi$ since (X_n) are i.i.d. By applying the above [note](#), we further have

$$\begin{aligned} (\phi(t/\sqrt{n}))^n &= \left(\phi(0) + \phi'(0) \frac{t}{\sqrt{n}} + \frac{t^2}{n} \int_0^1 \int_0^1 u\phi''(uvt/\sqrt{n}) du dv \right)^n \\ &= \left(1 + \frac{t^2}{n} \int_0^1 \int_0^1 u\phi''(uvt/\sqrt{n}) du dv \right)^n \end{aligned}$$

since $\phi(0) = 1$ and $\phi'(0) = i\mu = 0$. It remains to show that the double integral converges to $-1/2$ since it'll imply $(\phi(t/\sqrt{n}))^n \rightarrow e^{-t^2/2}$. We see that as $n \rightarrow \infty$, the integrand

$$u\phi''(uvt/\sqrt{n}) \rightarrow u\phi''(0) = u(i^2\mathbb{E}[X^2]) = -u(\text{Var}[X] + (\mathbb{E}[X])^2) = -u(1+0) = -u.$$

Hence, from the [bounded convergence theorem](#),

$$\int_0^1 \int_0^1 u\phi''(ut/\sqrt{n}) du dv \rightarrow \int_0^1 \int_0^1 -u du dv = -\frac{1}{2},$$

which shows the result. ■

Remark. From the [central limit theorem](#), we can indeed deduce the [weak law of large number](#). But since the former requires more conditions, hence [weak law of large number](#) still has its own merit.

2.4.1 Application to Inference

Firstly, let's consider some applications for mean estimation. Let X, X_1, \dots, X_n be i.i.d. samples such that $\mathbb{E}[X] = \mu$, $\text{Var}[X] = \sigma^2$. If, also, X_i 's are Gaussian, $\bar{X}_n \sim \mathcal{N}(\mu, \sigma^2/n)$, i.e.,

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1),$$

Intuition. We make the distribution independent of parameters to get a confidence interval.

However, the left-hand side is not an estimator now since it depends on σ . If we replace it by the sample standard deviation $\hat{\sigma}_n$, as $n \rightarrow \infty$,

$$T_n := \frac{\bar{X}_n - \mu}{\hat{\sigma}_n/\sqrt{n}} \sim t_{n-1} \xrightarrow{\text{TV}} \mathcal{N}(0, 1)$$

where T_n is the t -statistic.

Problem. What if X_i 's are not Gaussian?

Answer. We see that if $\hat{\sigma}_n^2 \xrightarrow{P} \sigma^2$, as $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} \mathcal{N}(0, 1)$ from [central limit theorem](#),

$$T_n = \frac{\sigma}{\hat{\sigma}_n} \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} \mathcal{N}(0, 1)$$

due to [Slutsky's theorem](#). Indeed, by letting $Y_i := X_i - \mu$ for all i (and also $Y = X - \mu$),

$$\frac{n-1}{n} \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - (\bar{Y}_n)^2$$

As $n \rightarrow \infty$, $(n-1)/n \rightarrow 1$, $\frac{1}{n} \sum_{i=1}^n Y_i^2 \xrightarrow{P} \mathbb{E}[Y^2] = \text{Var}[X] = \sigma^2$,^a and finally $(\bar{Y}_n)^2 \xrightarrow{P} (\mathbb{E}[Y])^2 = 0$, both from [weak law of large number](#). Hence, $\hat{\sigma}_n^2 \xrightarrow{P} \sigma^2$. ⊗

^aWe need to check whether the first moment exists.

Remark. For mean estimation, even if the data is not Gaussian, we're fine.

Next, let's consider variance estimation. Again, let X, X_1, \dots, X_n be i.i.d. Gaussian random samples,

$$(n-1) \frac{\hat{\sigma}_n^2}{\sigma^2} \stackrel{D}{=} \sum_{i=1}^{n-1} Z_i^2$$

where $(Z_{n-1}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. Firstly, since $\mathbb{E}[Z_i^2] = \text{Var}[Z_i] + (\mathbb{E}[Z_i])^2 = 1$, and $\text{Var}[Z_i^2] = \mathbb{E}[Z_i^4] - (\mathbb{E}[Z_i^2])^2 = 3 - 1 = 2$. Standardizing,

$$\frac{(n-1) \frac{\hat{\sigma}_n^2}{\sigma^2} - (n-1)}{\sqrt{2(n-1)}} \stackrel{D}{=} \frac{\sum_{i=1}^{n-1} Z_i^2 - (n-1)}{\sqrt{2(n-1)}} \xrightarrow{D} \mathcal{N}(0, 1),$$

i.e., as $n \rightarrow \infty$,

$$\sqrt{n-1} \left(\frac{\hat{\sigma}_n^2}{\sigma^2} - 1 \right) \xrightarrow{D} \mathcal{N}(0, 2) \Leftrightarrow \sqrt{n} \left(\frac{\hat{\sigma}_n^2}{\sigma^2} - 1 \right) \xrightarrow{D} \mathcal{N}(0, 2) \Leftrightarrow \sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \xrightarrow{D} \mathcal{N}(0, 2\sigma^4).$$

Let's postpone the problem of getting a confident interval, and first ask the following.

Problem. What if X_i 's are not Gaussian?

Answer. We see that from the same calculation as above, with $Y_i := X_i - \mu$ (and also $Y = X - \mu$),

$$\begin{aligned} \hat{\sigma}_n^2 - \frac{\hat{\sigma}_n^2}{n} &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2 \Rightarrow \hat{\sigma}_n^2 - \sigma^2 - \frac{\hat{\sigma}_n^2}{n} = \frac{1}{n} \sum_{i=1}^n (Y_i^2 - \sigma^2) - \bar{Y}_n^2 \\ &\Rightarrow \sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) - \frac{\hat{\sigma}_n^2}{\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^2 - \sigma^2) - \frac{(\sqrt{n}\bar{Y}_n)^2}{\sqrt{n}}. \end{aligned}$$

As $n \rightarrow \infty$, $\hat{\sigma}_n^2/\sqrt{n} \xrightarrow{P} 0$, and from the [central limit theorem](#), if $\mathbb{E}[X^4] < \infty$,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^2 - \sigma^2) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^2 - \mathbb{E}[Y_i^2]) \xrightarrow{D} \mathcal{N}(0, \text{Var}[Y_i^2]),$$

and finally, again from the [central limit theorem](#), $\sqrt{n}\bar{Y}_n$ [converges in distribution](#),

$$\frac{(\sqrt{n}\bar{Y}_n)^2}{\sqrt{n}} = o(1)O_p(1) = o_p(1).$$

Hence, as long as $\mathbb{E}[X^4] < \infty$, then

$$\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \xrightarrow{D} \mathcal{N}(0, \text{Var}[Y_i^2]).$$

We see that

$$\text{Var}[Y^2] = \mathbb{E}[(X - \mu)^4] - (\mathbb{E}[(X - \mu)^2])^2 = \sigma^4 \mathbb{E} \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] - \sigma^4 = \sigma^4 \left(\mathbb{E} \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] - 1 \right),$$

so we need the *Kurtosis* of X , i.e., $\mathbb{E}[(X - \mu)/\sigma]^4$. Previously, when X_i 's are Gaussian, it is 3. \circledast

We note the following.

Note. Let $Z = (X - \mu)/\sigma$, then the above is meaningful if $\mathbb{E}[Z^4] > 1$ since otherwise the variance is 0. However, from Jensen's inequality, $\mathbb{E}[Z^4] \geq (\mathbb{E}[Z^2])^2 \geq 1$, hence the above makes sense when the equality doesn't hold.

Example. The equality holds (i.e., $\mathbb{E}[Z^4] = 1$) when $Z^2 = 1$, or $Z \in \{\pm 1\}$. Hence, the above might happen for Z being binary.

The takeaway is that if the Kurtosis of X is different from the normal, then the distribution of $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2)$ is different. Moreover, we see that we can't quite do the same trick as before by making the distribution free from parameters to get a confident interval.

Note. If we don't know the Kurtosis of X , we can't say anything about the confident interval.

Lecture 11: Testing Normality

Let's try to generalize what we have done. Let $Y := X - \mu = X - \mathbb{E}[X]$ (and also $Y_i = X_i - \mu$ as usual), $\mu_k := \mathbb{E}[Y^k] = \mathbb{E}[(X - \mu)^k]$ for all $k \geq 2$, and finally $\tilde{\mu}_k = \mu_k / \sigma^k = \mathbb{E}[(X - \mu)^k / \sigma^k]$. 20 Feb. 9:30

As previously seen. In this notation, we proved $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \rightarrow \mathcal{N}(0, (\tilde{\mu}_4 - 1)\sigma^4)$, or equivalently,

$$\frac{\sqrt{n}}{\sqrt{\tilde{\mu}_4 - 1}} \left(\frac{\hat{\sigma}_n^2}{\sigma^2} - 1 \right) \rightarrow \mathcal{N}(0, 1).$$

As hinted before, we now address the issue of independence of parameters on the left-hand side. Since we know how to estimate σ (i.e., by $\hat{\sigma}_n$) **consistently**, it reduces to estimating μ_k , so we ask the following.

Problem. How to estimate $\tilde{\mu}_4$, or more generally, how to estimate $\tilde{\mu}_k$ **consistently**?

Consider the k^{th} sample centered moment

$$M_k := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^k.$$

Intuition. If we know μ , then $\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^k \xrightarrow{P} \mu_k$ by the **weak law of large number**.

However, since we don't know μ , we need to use \bar{X}_n . But this still yields a **consistent** estimator.

Proposition 2.4.1. If $\mu_k < \infty$, $M_k \xrightarrow{P} \mathbb{E}[Y^k] = \mu_k$.

Proof. Let's denote $\bar{X}_n =: \bar{X}$ and $\bar{Y}_n =: \bar{Y}$. Then

$$M_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^k = \frac{1}{n} \sum_{i=1}^n \sum_{\ell=0}^k \binom{k}{\ell} Y_i^\ell (-\bar{Y})^{k-\ell} = \sum_{\ell=0}^k \binom{k}{\ell} (-\bar{Y})^{k-\ell} \frac{1}{n} \sum_{i=1}^n Y_i^\ell.$$

Let $\frac{1}{n} \sum_{i=1}^n Y_i^\ell =: \bar{Y}^\ell$, then we further get

$$M_k = \sum_{\ell=0}^k \binom{k}{\ell} (-\bar{Y})^{k-\ell} \bar{Y}^\ell = \bar{Y}^k + \sum_{\ell=0}^{k-1} \binom{k}{\ell} (-\bar{Y})^{k-\ell} \bar{Y}^\ell.$$

By the **weak law of large number**, $\bar{Y}^k \xrightarrow{P} \mathbb{E}[Y^k] = \mu_k$ and $(-\bar{Y})^{k-\ell} \xrightarrow{P} 0$ since $-\bar{Y} \xrightarrow{P} 0$ with **continuous mapping theorem**, hence $M_k \xrightarrow{P} \mu_k$ by **Slutsky's theorem**. ■

Furthermore, we can also ask for the asymptotic distribution of M_k , i.e., $\sqrt{n}(M_k - \mu_k)$. Firstly,

$$\sqrt{n}(M_k - \mu_k) = \sqrt{n}(\bar{Y}^k - \mu_k) + \sum_{\ell=0}^{k-1} \binom{k}{\ell} (-\bar{Y})^{k-\ell} \bar{Y}^\ell \sqrt{n} = \sqrt{n}(\bar{Y}^k - \mu_k) + \sum_{\ell=0}^{k-1} \binom{k}{\ell} \frac{(-\bar{Y} \sqrt{n})^{k-\ell}}{\sqrt{n}^{k-\ell-1}} \bar{Y}^\ell.$$

We see that

- $\bar{Y}\sqrt{n}$ is asymptotically normal, hence is [bounded in probability](#) from [Proposition 2.2.5](#);
- $\bar{Y}^\ell \xrightarrow{p} \mathbb{E}[Y^\ell] = O(1)$;
- $1/\sqrt{n}^{k-\ell-1} = o(1)$ for $\ell < k-1$.

Combining, every term in the summation is $O(1)O_p(1)o(1) = o_p(1)$ except for $\ell = k-1$. This means

$$\begin{aligned}\sqrt{n}(M_k - \mu_k) &= \sqrt{n}(\bar{Y}^k - \mu_k) - \binom{k}{k-1} \bar{Y}^{k-1} \sqrt{n}\bar{Y} + \sum_{\ell=0}^{k-2} \binom{k}{\ell} o_p(1) \\ &= \sqrt{n}(\bar{Y}^k - \mu_k) - k \bar{Y}^{k-1} \sqrt{n}\bar{Y} + o_p(1)\end{aligned}$$

The problem is that although $\sqrt{n}\bar{Y} = O_p(1)$, \bar{Y}^{k-1} is not $o_p(1)$. By replacing \bar{Y}^{k-1} by $\bar{Y}^{k-1} - \mu_{k-1} + \mu_{k-1}$,

$$\begin{aligned}&= \sqrt{n}(\bar{Y}^k - \mu_k) - k \left(\bar{Y}^{k-1} - \mu_{k-1} \right) \sqrt{n}\bar{Y} - k\mu_{k-1} \sqrt{n}\bar{Y} + o_p(1) \\ &= \sqrt{n}(\bar{Y}^k - \mu_k) - k\mu_{k-1} \sqrt{n}\bar{Y} + o_p(1)\end{aligned}$$

since $\bar{Y}^{k-1} - \mu_{k-1} \xrightarrow{p} 0$ from the [weak law of large number](#), finally,

$$\begin{aligned}&= \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n (Y_i^k - \mu_k) \right) - k\mu_{k-1} \frac{1}{n} \sqrt{n} \sum_{i=1}^n Y_i + o_p(1) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^k - \mu_k - k\mu_{k-1} Y_i) + o_p(1).\end{aligned}$$

Observe that $Y_i^k - \mu_k - k\mu_{k-1} Y_i$'s are i.i.d., the whole thing converges to $\mathcal{N}(0, \text{Var}[Y^k - \mu_k - k\mu_{k-1} Y])$ [in distribution](#) by [central limit theorem](#), where the variance can be calculated as

$$\begin{aligned}\text{Var}[Y^k - \mu_k - k\mu_{k-1} Y] &= \text{Var}[Y^k - k\mu_{k-1} Y] \\ &= \text{Var}[Y^k] + k^2 \mu_{k-1}^2 \text{Var}[Y] - 2k\mu_{k-1} \text{Cov}[Y, Y^k] \\ &= \mu_{2k} - \mu_k^2 + k^2 \mu_{k-1}^2 \sigma^2 - 2k\mu_{k-1} \mu_{k+1}\end{aligned}$$

since $\mathbb{E}[Y] = 0$, $\text{Var}[Y] = \sigma^2$, and $\text{Cov}[Y, Y^k] = \mathbb{E}[Y \cdot Y^k] - \mathbb{E}[Y]\mathbb{E}[Y^k] = \mathbb{E}[Y^{k+1}] = \mu_{k+1}$.

Remark. We didn't give a confidence interval for estimating σ^2 , which requires we to obtain the joint distribution of $\hat{\sigma}_n^2$ and M_k .

2.4.2 Testing Normality with Odd Moments

Instead of giving a confidence interval for $\hat{\sigma}_n^2$, let's consider using two distributional results for M_k and $\hat{\sigma}_n^2$ and apply it to the problem of testing normality, i.e., let $H_0: X \sim \mathcal{N}(\mu, \sigma^2)$ for some μ, σ^2 .

Intuition. With [Proposition 2.4.1](#), the idea is that to reject H_0 if $|M_k|$ is "large".

However, it turns out that considering $|M_k/\hat{\sigma}_n^k|$ is more appropriate. Hence, we again need to compute the joint of M_k and $\hat{\sigma}_n^k$, so we won't miss anything indeed. Anyway, now the problem is the following.

Problem. What is the asymptotic distribution of $M_k/\hat{\sigma}_n^k$?

First observe that since Gaussian is symmetric, $\mu_k = 0$ (and hence $\tilde{\mu}_k = 0$) for all odd k . It turns out that this property allows us to bypass the joint if we focus on odd k .

As previously seen. Previously we have $\sqrt{n}(M_k - \mu_k) \xrightarrow{D} \mathcal{N}(0, \text{Var}[Y^k - k\mu_{k-1} Y])$.

Formally, suppose k is odd, and $X \sim \mathcal{N}(\mu, \sigma^2)$, then $\mu_k = 0$, hence

$$\sqrt{n}(M_k - \mu_k) \xrightarrow{D} \mathcal{N}(0, \text{Var}[Y^k - k\mu_{k-1} Y]) \Rightarrow \sqrt{n} \frac{M_k}{\sigma^k} \xrightarrow{D} \mathcal{N}(0, \sigma^{-2k} \text{Var}[Y^k - k\mu_{k-1} Y])$$

By [Slutsky's theorem](#), $\sqrt{n}M_k/\hat{\sigma}^k$ also [converges](#) to this normal. Now, in particular, if $k = 3$, we have

$$\sqrt{n} \frac{M_3}{\hat{\sigma}_n^3} \xrightarrow{D} \mathcal{N}(0, \sigma^{-6} \text{Var}[Y^3 - 3\sigma^2 Y]) = \mathcal{N}(0, \sigma^{-6} (\text{Var}[Y^3] + 9\sigma^4 \sigma^2 - 6\sigma^2 \mathbb{E}[Y^4]))$$

where $\mu_2 = \sigma^2$ and $\text{Cov}[Y^3, Y] = \mathbb{E}[Y^4] - \mathbb{E}[Y]\mathbb{E}[Y^3] = \mathbb{E}[Y^4]$. We note the following.

Note. For odd k , $\text{Var}[Y^k] = \mathbb{E}[Y^{2k}] - (\mathbb{E}[Y^k])^2 = \mathbb{E}[Y^{2k}] = \mu_{2k}$ since $(\mathbb{E}[Y^k])^2 = \mu_k^2 = 0$.

Hence, by plugging $\text{Var}[Y^3] = \mu_{2 \times 3} = \mu_6$, the variance of the normal is further equal to

$$\frac{\mu_6 + 9\sigma^6 - 6\sigma^2 \mu_4}{\sigma^6} = \tilde{\mu}_6 + 9 - 6\tilde{\mu}_4 = 15 + 9 - 6 \times 3 = 6,$$

i.e.,

$$\sqrt{\frac{n}{6}} \frac{M_3}{\hat{\sigma}_n^3} \xrightarrow{D} \mathcal{N}(0, 1).$$

Remark. We get the asymptotic distribution of $M_k/\hat{\sigma}_n^k$ without computing the joint of M_k and $\hat{\sigma}_n^k$.

For even k , we really need to work out the joint. Since we know the asymptotic distribution of both M_k and $\hat{\sigma}_k$, the joint with can be obtained by $g(M_k, \hat{\sigma}_2) = |M_k/\hat{\sigma}^k|$ with the [delta method](#), with the multivariate version of [Theorem 2.4.2](#) since we now have two quantities.

2.4.3 Multivariate Central Limit Theorem

Our next goal is to prove the [multivariate central limit theorem](#), i.e., the high dimensional generalization of [Theorem 2.4.2](#). We first need the following tool.

Theorem 2.4.3 (Cramér-Wold device). Let (X_n) be a sequence of random vectors and X be a random vector in \mathbb{R}^d . Then $X_n \xrightarrow{D} X$ if and only if $t \cdot X_n \xrightarrow{D} t \cdot X$ for every $t \in \mathbb{R}^d$.

Proof. The forward direction is clear from [continuous mapping theorem](#) for the linear functional induced from t . For the backward direction, assume that $t \cdot X_n \xrightarrow{D} t \cdot X$. Then

$$\phi_{X_n}(t) = \mathbb{E}[e^{it \cdot X_n}] = \phi_{t \cdot X_n}(1) \rightarrow \phi_{t \cdot X}(1) = \mathbb{E}[e^{it \cdot X}] = \phi_X(t),$$

which implies $X_n \xrightarrow{D} X$ by the [Lévy-Cramer continuity theorem](#). ■

Remark. Proving $X_n \xrightarrow{D} X$ reduces to proving something in the scalar case.

Theorem 2.4.4 (Multivariate central limit theorem). Let (X_n) be i.i.d. random vectors in \mathbb{R}^d with $\mathbb{E}[X_i] = \mu \in \mathbb{R}^d$, $\text{Var}[X_i] = \Sigma \in \mathbb{R}^{d \times d}$ for all $1 \leq i \leq n$. Then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \xrightarrow{D} \mathcal{N}(0, \Sigma).$$

Lecture 12: Asymptotic Joint Distribution by Multivariate CLT

22 Feb. 9:30

Proof. Set $\mu = 0$, and it suffices to show that for any $t \in \mathbb{R}^d$,

$$t \cdot \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \right) \xrightarrow{D} t \cdot Z \sim \mathcal{N}(0, t^\top \Sigma t)$$

where $Z \sim \mathcal{N}(0, \Sigma)$. We see that from the [univariate central limit theorem](#), the left-hand side is

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n t \cdot X_i \xrightarrow{D} \mathcal{N}(0, \text{Var}[t \cdot X_i]),$$

and since $\text{Var}[t \cdot X] = t^\top \text{Var}[X] t = t^\top \Sigma t = \text{Var}[t \cdot Z]$, hence we're done. \blacksquare

2.4.4 Testing Normality with General Moments

With [multivariate central limit theorem](#), we can continue on the problem of computing the asymptotic distribution of $\widetilde{M}_k := M_k / \hat{\sigma}_n^k$ for general k . Recall the setup, where we let (X_n) and X be i.i.d. random variable, $Y_i = X_i - \mu$ (and $Y = X - \mu$), $\sigma^2 = \text{Var}[X]$, $\mu_k = \mathbb{E}[Y^k]$, and $\tilde{\mu}_k = \mu_k / \sigma^k$. Let's start with $k = 1$, i.e., compute the asymptotic law of $\bar{X}_n / \hat{\sigma}_n$.

As previously seen. We have proved that

- $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \mathcal{N}(0, \sigma^2)$ from $\sqrt{n}(\bar{X}_n - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$;
- $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \xrightarrow{D} \mathcal{N}(0, \mu_4 - \sigma^4)$ from $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^2 - \sigma^2) + o_p(1)$.^a

^aAssuming that $\mu_4 - \sigma^4 > 1$.

We see that we have \bar{X}_n and $\hat{\sigma}_n^2$, not $\hat{\sigma}_n$. But this is fine since we can

1. compute the joint distribution of \bar{X}_n and $\hat{\sigma}_n^2$;
2. apply the [delta method](#) with $g(\bar{X}_n, \hat{\sigma}_n^2) := \bar{X}_n / \hat{\sigma}_n$ to get the distribution of $\bar{X}_n / \hat{\sigma}_n$.

Specifically, from the [multivariate central limit theorem](#) and above result we proved,

$$\sqrt{n} \left(\begin{pmatrix} \bar{X}_n \\ \hat{\sigma}_n^2 \end{pmatrix} - \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \right) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} Y_i \\ Y_i^2 - \sigma^2 \end{pmatrix} + o_p(1) \xrightarrow{D} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \text{Var} \left[\begin{pmatrix} Y \\ Y^2 \end{pmatrix} \right] \right),$$

such that

$$\text{Var} \left[\begin{pmatrix} Y \\ Y^2 \end{pmatrix} \right] = \begin{pmatrix} \text{Var}[Y] & \text{Cov}[Y, Y^2] \\ \text{Cov}[Y, Y^2] & \text{Var}[Y^2] \end{pmatrix} = \begin{pmatrix} \sigma^2 & \mu_3 \\ \mu_3 & \mu_4 - \sigma^4 \end{pmatrix}.$$

Let's leave the application of the [delta method](#) to the general k .

Intuition. The actual characterization of \bar{X}_n and $\hat{\sigma}_n^2$ right before applying [central limit theorem](#) is much more useful than the final asymptotic distributions.

Next, we compute the asymptotic law of $\widetilde{M}_k = M_k / \hat{\sigma}_k$, where we have already proven the following.

As previously seen. $\sqrt{n}(M_k - \mu_k) \rightarrow \mathcal{N}(0, \text{Var}[Y^k - k\mu_{k-1}Y])$ from

$$\sqrt{n}(M_k - \mu_k) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i^k - \mu_k - k\mu_{k-1}Y_i) + o_p(1).$$

Then as above, we again use the result for $\hat{\sigma}_n^2$ and get

$$\begin{aligned} Y := \sqrt{n} \left(\begin{pmatrix} \hat{\sigma}_n^2 \\ M_k \end{pmatrix} - \begin{pmatrix} \sigma^2 \\ \mu_k \end{pmatrix} \right) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} Y_i - \sigma^2 \\ Y_i^k - \mu_k - k\mu_{k-1}Y_i \end{pmatrix} + o_p(1) \\ &\xrightarrow{D} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \text{Var}[Y^2] & \text{Cov}[Y^2, Y^k - k\mu_{k-1}Y] \\ \text{Cov}[Y^2, Y^k - k\mu_{k-1}Y] & \text{Var}[Y^k - k\mu_{k-1}Y] \end{pmatrix} \right). \end{aligned}$$

This “Y” will be used in the [delta method](#) later.²

²This is not exact since Y should be the random vector corresponding the asymptotic distribution on the right-hand side. But this is fine in the end as we will soon see.

Remark. If $\mu_\ell = 0$ for all odd ℓ , then M_k and $\hat{\sigma}_n^2$ are “asymptotically independent”. This is why we get a simplification for odd case.

Proof. Since then $\text{Cov}[Y^2, Y^k - k\mu_{k-1}Y] = 0$ for odd k . *

Now, since we’re using $\hat{\sigma}_n^2$ but not $\hat{\sigma}_n^k$, we need to use [delta method](#) by considering

$$\widetilde{M}_k = \frac{M_k}{\hat{\sigma}_n^k} =: g(\hat{\sigma}_n^2, M_k)$$

where $g(x, y) = y/x^{k/2}$ for $x > 0, y \in \mathbb{R}$. We see that

$$\nabla g(\sigma^2, \mu_k) = \begin{pmatrix} -\frac{k}{2}\mu_k\sigma^{-k-2} & \sigma^{-k} \end{pmatrix} = \begin{pmatrix} -\frac{k}{2}\tilde{\mu}_k\sigma^{-2} & \sigma^{-k} \end{pmatrix}$$

since $\tilde{\mu}_k = \mu_k/\sigma^k$ and

- $\partial g/\partial x = -k y x^{-k/2-1}/2$.
- $\partial g/\partial y = x^{-k/2}$.

From [delta method](#), with $\tilde{\mu}_k = g(\sigma^2, \mu_k)$, we get $\sqrt{n}(g(\hat{\sigma}_n^2, M_k) - g(\sigma^2, \mu_k)) \xrightarrow{D} \nabla g Y$, or equivalently,

$$\begin{aligned} \sqrt{n}(\widetilde{M}_k - \tilde{\mu}_k) &= \nabla g(\sigma^2, \mu_k) \frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} Y_i^k - \mu_k - k\mu_{k-1}Y_i \\ Y_i^2 - \sigma^2 \end{pmatrix} + o_p(1) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(-\frac{k}{2}\tilde{\mu}_k \frac{1}{\sigma^2} (Y_i^2 - \sigma^2) + \frac{1}{\sigma^k} (Y_i^k - \mu_k - k\mu_{k-1}Y_i) \right) + o_p(1) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(-\frac{k}{2}\tilde{\mu}_k (Z_i^2 - 1) + (Z_i^k - \tilde{\mu}_k - k\tilde{\mu}_{k-1}Z_i) \right) + o_p(1) \xrightarrow{D} \mathcal{N}(0, \theta_k) \end{aligned}$$

by letting $Z_i := (X_i - \mu)/\sigma = Y_i/\sigma$ and

$$\theta_k := \text{Var} \left[-\frac{k}{2}\tilde{\mu}_k (Z^2 - 1) + (Z^k - \tilde{\mu}_k - k\tilde{\mu}_{k-1}Z) \right],$$

which is independent of parameters since under H_0 , $Z \sim \mathcal{N}(0, 1)$.

Note. It’s more convenient to use [delta method](#) in this way, i.e., not use the actual Y corresponding to the distribution, but use the term in the limiting sequence with $o_p(1)$. Compared to the last time, we do this for odd k , and only get the asymptotic distribution, not this explicit decomposition.

With this explicit formula, we now see one example.

Example. Consider using both \widetilde{M}_3 and \widetilde{M}_4 to test $H_0: X \sim \mathcal{N}$. We see that under H_0 ,

$$\left(\sqrt{\frac{n}{\theta_3}} \widetilde{M}_3 \right)^2 + \left(\sqrt{\frac{n}{\theta_4}} (\widetilde{M}_4 - \mu_4) \right)^2 \xrightarrow{D} \chi_2^2.$$

Proof. One can write down $\sqrt{n}(\widetilde{M}_\ell, \tilde{\mu}_\ell)$ for even ℓ , and also $\sqrt{n}(\widetilde{M}_k - \tilde{\mu}_k)$ for odd k , and see that the covariance indeed is 0, hence independent, so they add up to χ_2^2 . *

2.4.5 Asymptotic Relative Efficiency

Assume $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Pois}(\theta)$. To estimate θ , as $\theta = \mathbb{E}[X] = \text{Var}[X]$, two natural estimators are \bar{X}_n and $\hat{\sigma}_n^2$. To compare them, we see that

- $\sqrt{n}(\bar{X}_n - \theta) \xrightarrow{D} \mathcal{N}(0, \sigma^2)$;
- $\sqrt{n}(\hat{\sigma}_n^2 - \theta) \xrightarrow{D} \mathcal{N}(0, \mu_4 - \sigma^4)$.

As $\sigma^2 = \theta$ and $\mu_4 = 3\theta^2 + \theta$, we see that \bar{X}_n is better since its variance is smaller.

Problem. But by how much?

Answer. We can consider how many data we need such that we get a similar precision. Consider

$$\sqrt{n}(T_n^i - \theta) \xrightarrow{D} \mathcal{N}(0, \sigma_i^2(\theta))$$

for two statistics T_i , $i = 1, 2$. This implies

$$\mathbb{P}\left(\theta \in T_n^i \pm Z_{\alpha/2} \frac{\sigma_i(\theta)}{\sqrt{n}}\right) \cong 1 - \alpha.$$

Let $I_n^{(i)} := T_n^i \pm Z_{\alpha/2} \sigma_i(\theta)/\sqrt{n}$, and let n_i be the value of n such that $|I_n^{(i)}| = \gamma$,

$$\gamma = 2Z_{\alpha/2} \frac{\sigma_i(\theta)}{\sqrt{n_i}} \Rightarrow n_i = \left(\frac{2Z_{\alpha/2} \sigma_i(\theta)}{\gamma}\right)^2,$$

i.e., $n_1/n_2 = \sigma_1(\theta)^2/\sigma_2(\theta)^2$, which we called the *asymptotic relative efficiency* $\text{ARE}(T^1, T^2)$. *

Now, we consider \bar{X}_n as the estimator of θ . We have $\sqrt{n}(\bar{X}_n - \theta) \xrightarrow{D} \sqrt{\theta}\mathcal{N}(0, 1) = \mathcal{N}(0, \theta)$.

Note. As the asymptotic distribution depends on θ , we don't directly get a confidence interval.

To get around this, we first write

$$\sqrt{n}(\bar{X}_n - \theta) \xrightarrow{D} \sqrt{\theta}Z \sim \mathcal{N}(0, \theta)$$

for $Z \sim \mathcal{N}(0, 1)$. Then, instead of writing this as

$$\frac{\sqrt{n}}{\sqrt{\theta}}(\bar{X}_n - \theta) \xrightarrow{D} Z$$

and replace $\sqrt{\theta}$ on by some estimator, observe that from [delta method](#) with some g ,

$$\sqrt{n}(g(\bar{X}_n) - g(\theta)) \xrightarrow{D} g'(\theta)\sqrt{\theta}Z$$

Suppose $g'(\theta)\sqrt{\theta} = c$ is some constant for every $\theta > 0$, our goal is also achieved, i.e.,

$$\frac{\sqrt{n}}{c}(g(\bar{X}_n) - g(\theta)) \xrightarrow{D} \mathcal{N}(0, 1).$$

Claim. There exists g such that $c = 1/2$.

Proof. Since for $g'(\theta) = \frac{1}{2\sqrt{\theta}}$, we have $g(\theta) = \sqrt{\theta}$. *

Then, the asymptotic confidence interval for $g(\theta)$ is just

$$\left(g(\bar{X}_n) - Z_{\alpha/2} \frac{c}{\sqrt{n}}, g(\bar{X}_n) + Z_{\alpha/2} \frac{c}{\sqrt{n}}\right),$$

and apply g^{-1} to get back θ . In this case, $g^{-1}(u) = u^2$. This is the so-called *variance stabilizing transformation*.

Appendix

Bibliography

- [Das08] Anirban DasGupta. *Asymptotic Theory of Statistics and Probability*. Springer Science & Business Media, Feb. 6, 2008. 727 pp. ISBN: 978-0-387-75971-5. Google Books: [sX4_AAAAQBAJ](#).
- [Fer17] Thomas S. Ferguson. *A Course in Large Sample Theory*. Routledge, Sept. 6, 2017. 140 pp. ISBN: 978-1-351-47005-6. Google Books: [clcODwAAQBAJ](#).
- [Leh04] E. L. Lehmann. *Elements of Large-Sample Theory*. Springer Science & Business Media, Aug. 27, 2004. 640 pp. ISBN: 978-0-387-98595-4. Google Books: [geIoxvgTXlEC](#).
- [Ser09] Robert J. Serfling. *Approximation Theorems of Mathematical Statistics*. John Wiley & Sons, Sept. 25, 2009. 399 pp. ISBN: 978-0-470-31719-8. Google Books: [enUouJ4EHzQC](#).
- [Vaa98] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press, 1998. ISBN: 978-0-521-78450-4. DOI: [10.1017/CB09780511802256](#). URL: <https://www.cambridge.org/core/books/asymptotic-statistics/A3C7DAD3F7E66A1FA60E9C8FE132EE1D> (visited on 10/17/2023).