

Review of Statistics

Part A: Properties of Estimators

THOMAS WIEMANN
University of Chicago

Econometrics
Econ 21020

Updated: April 13, 2022

Recap

The review of probability theory introduced a formal language for characterizing uncertainty.

- ▷ Introduced random variables and their probability distributions;
- ▷ Developed concepts to describe features of random variables;
- ▷ Discussed restrictions on the joint distribution of random variables.

With our toolbox, we can return to the returns to education example.

- ▷ Under the random assignment assumption, we can show that

$$E_U[g(1, U) - g(0, U)|W = 1] = E_Y[Y|W = 1] - E_Y[Y|W = 0],$$

where $E_Y[Y|W = 1]$ and $E_Y[Y|W = 0]$ are features of the joint distribution of the observables (Y, W) .

Note that $E_Y[Y|W = 1]$ and $E_Y[Y|W = 0]$ are *theoretical* concepts.

- ▷ Statistics forms a bridge between random variables and data.

1. Estimators
2. Finite Sample Properties
 - ▷ Bias
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

These notes benefit greatly from the exposition in Wasserman (2003) and the lecture notes of Prof. Max Tabord-Meehan.

1. **Estimators**
2. Finite Sample Properties
 - ▷ Bias
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

These notes benefit greatly from the exposition in Wasserman (2003) and the lecture notes of Prof. Max Tabord-Meehan.

Random Sampling

Consider independent random variable X_1, \dots, X_n with $X_i \sim F_i, \forall i$.

- ▷ When $F_i = F, \forall i = 1, \dots, n$, we say that X_1, \dots, X_n are *independent and identically distributed* (iid).
- ▷ To denote an iid sample of size n from F , we write

$$X_1, \dots, X_n \stackrel{iid}{\sim} F. \quad (1)$$

Example 1

Consider $X_1 \sim N(\mu_1, \sigma_1^2)$ and $X_2 \sim N(\mu_2, \sigma_2^2)$.

- ▷ If $X_1 \perp X_2$, then X_1 and X_2 are independent.
- ▷ If $(\mu_1, \sigma_1^2) = (\mu_2, \sigma_2^2)$, then X_1 and X_2 are identically distributed.
- ▷ If $X_1 \perp X_2$ and $(\mu_1, \sigma_1^2) = (\mu_2, \sigma_2^2)$, then X_1 and X_2 are iid.

Notation: Instead of (1), we also sometimes write $X_1, \dots, X_n \stackrel{iid}{\sim} X$. So X may denote a random variable or its distribution.

Estimators

Statistics is concerned with learning about the distribution from F using a sample $X_1, \dots, X_n \sim F$.

- ▷ We will (for the most part), consider iid-samples.

Instead of fully characterizing F , the focus often lies on features of F .

- ▷ Features of interest are called *parameters*.
- ▷ For example, we may be interested in $\mu \equiv E[X]$ where $X \sim F$. Here, μ is the parameter of interest.

An *estimate* is a “guess” for the value of the parameter of interest.

- ▷ An *estimator* is a function of the sample whose value serves as a “guess” for a parameter of interest.
- ▷ For example, if $\mu \in \mathbb{R}$ and $\text{supp } X_i = \mathbb{R}, \forall i$, then an estimator for μ is a function $\hat{\mu}_n(X_1, \dots, X_n)$.
- ▷ Importantly: μ is a number but $\hat{\mu}_n$ is a random variable.

Notation: *Subscripts on expectation operators or distribution functions are omitted from now on whenever the context is clear.*

Example 2

Consider a sample $X_1, \dots, X_n \stackrel{iid}{\sim} F$. An estimator for $F(x) = P(X \leq x)$ is given by

$$\widehat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}, \quad (2)$$

that is, the share of the sample below x is a “guess” for $P(X \leq x)$.

The estimator \widehat{F}_n is called the *empirical CDF*.

The empirical CDF leads to a class of estimators that are known under the *sample analogue principle*.

- ▷ Suppose we are interested in a feature of F . The sample analogue principle suggests using the analogous feature of \widehat{F}_n as an estimate.

Example 3

Consider a sample $X_1, \dots, X_n \stackrel{iid}{\sim} F$. Let $\mu = E[X]$ denote the parameter of interest. The sample analogue principle suggests the estimator

$$\hat{\mu}_n \equiv E_n[X] = \frac{1}{n} \sum_{i=1}^n X_i, \quad (3)$$

where E_n denotes the expectation with respect to the empirical CDF \hat{F}_n .

Similarly, if the parameter of interest is $\sigma^2 = \text{Var}(X)$, the sample analogue principle suggests the estimator

$$\hat{\sigma}_n^2 \equiv \quad (4)$$

Estimators (Contd.)

The sample analogue principle is not the only approach to constructing estimators. Another frequently encountered class of estimators are extremum estimators, defined as the minimizers of a loss-functions.

Example 4

Consider a sample $X_1, \dots, X_n \stackrel{iid}{\sim} F$ and let $\mu = E[X]$ denote the parameter of interest. Define an estimator

$$\hat{\mu}_n = \arg \min_{\mu \in \mathbb{R}} \sum_{i=1}^n (X_i - \mu)^2. \quad (5)$$

Taking first order conditions, we have

$$0 =$$

Estimators (Contd.)

For a given parameter, there infinitely many possible estimators.

Example 5

Consider a sample $X_1, \dots, X_n \stackrel{iid}{\sim} F$ and let $\mu = E[X]$ denote the parameter of interest. Each of the following are estimators for μ :

- ▷ $\hat{\mu}_n^{(1)} = 0$;
- ▷ $\hat{\mu}_n^{(2)} = X_1$;
- ▷ $\hat{\mu}_n^{(3)} = \frac{1}{n} \sum_{i=1}^n X_i$.
- ▷ $\hat{\mu}_n^{(4)} = \frac{1}{n+\lambda} \sum_{i=1}^n X_i$ for some fixed $\lambda \in \mathbb{R}_+$.

Which one do you like best?

Statistics provides tools that allow for comparisons of estimators.

- ▷ Allows for selecting the “best” (or – at least – a “good enough”) estimator.

Sampling Distribution

Recall that an estimator is a function of random variables and hence itself a random variable.

- ▷ The *sampling distribution* of an estimator is a name for its distribution.

Comparisons of estimators are analogous to comparisons of (features of) their sampling distribution.

- ▷ The sampling distribution often depends on the sample size n .

Consider an estimator $\hat{\theta}_n$ for some parameter θ of a distribution F .

- ▷ *Finite sample properties* describe features of the distribution of $\hat{\theta}_n$. These properties hold for any sample size $n \in \mathbb{N}$.
- ▷ *Large sample properties* describe features of the *asymptotic* distribution of $\hat{\theta}_n$. These properties hold approximately for large enough sample sizes n .

1. Estimators
2. **Finite Sample Properties**
 - ▷ **Bias**
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

We begin with describing the expected deviations of the estimator from the true parameter.

Definition 1

The *bias* of an estimator $\hat{\theta}_n$ for θ is defined as

$$\text{Bias}(\hat{\theta}_n) = E[\hat{\theta}_n] - \theta. \quad (6)$$

The estimator is said to be

- ▷ *unbiased* if $\text{Bias}(\hat{\theta}_n) = 0$;
- ▷ *downwards biased* if $\text{Bias}(\hat{\theta}_n) < 0$;
- ▷ *upwards biased* if $\text{Bias}(\hat{\theta}_n) > 0$.

Example 6

Consider the estimators $\hat{\mu}_n^{(1)}$, $\hat{\mu}_n^{(2)}$, $\hat{\mu}_n^{(3)}$ and $\hat{\mu}_n^{(4)}$ of Example 5. We have

$$\text{Bias}(\hat{\mu}_n^{(1)}) =$$

$$\text{Bias}(\hat{\mu}_n^{(2)}) =$$

$$\text{Bias}(\hat{\mu}_n^{(3)}) =$$

$$\text{Bias}(\hat{\mu}_n^{(4)}) =$$

Note that the Bias of $\hat{\mu}_n^{(4)}$ depends on the unknown parameter μ .

Example 7

Consider the estimator $\hat{\sigma}_n^2$ defined in Example 3. We have

$$\hat{\sigma}_n^2 =$$

and

$$\text{Bias}(\hat{\sigma}_n^2) =$$

Can you construct an unbiased estimate for $\text{Var}(X)$?

1. Estimators
2. **Finite Sample Properties**
 - ▷ Bias
 - ▷ **Variance**
 - ▷ The Bias-Variance Trade-Off
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

Example 6 showed that very different estimators can have the same bias.

- ▷ Require other features of the sampling distribution to make comparison useful.

Another key property of an estimator is its variance:

$$\text{Var}(\hat{\theta}_n) = E\left[\left(\hat{\theta}_n - E[\hat{\theta}_n]\right)^2\right] \quad (7)$$

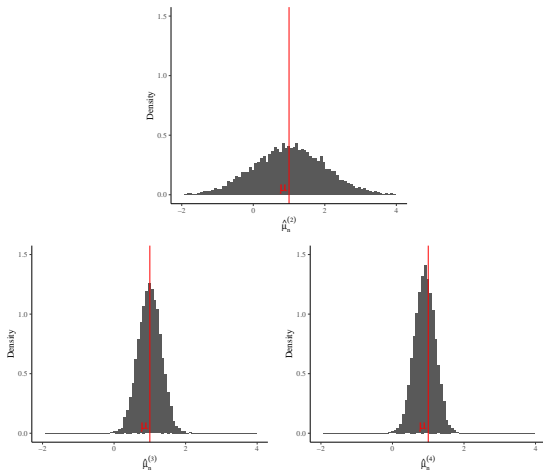
- ▷ Describes deviations from the expected value of the estimator.
- ▷ The expected value of a biased estimator is *not* the true parameter.

Figure 1 illustrates why considering both bias and variance is useful for distinguishing estimators.

- ▷ Draws from the sampling distribution of the estimators of Example 5.

Estimation Variance (Contd.)

Figure 1: Draws from Sampling Distributions of Estimators



Notes. Histograms of $\hat{\mu}_n^{(2)}$, $\hat{\mu}_n^{(3)}$ and $\hat{\mu}_n^{(4)}$ of Example 5 where $n = 10$ and $(\mu, \sigma^2) = (1, 1)$. For $\hat{\mu}_n^{(4)}$, I set $\lambda = 1$. You can find the corresponding code on GitHub: [lecture_plots.R](#).

Example 8

Consider the estimators $\hat{\mu}_n^{(1)}$, $\hat{\mu}_n^{(2)}$, $\hat{\mu}_n^{(3)}$ and $\hat{\mu}_n^{(4)}$ of Example 5. We have

$$\text{Var}(\hat{\mu}_n^{(1)}) =$$

$$\text{Var}(\hat{\mu}_n^{(2)}) =$$

$$\text{Var}(\hat{\mu}_n^{(3)}) =$$

$$\text{Var}(\hat{\mu}_n^{(4)}) =$$

Note that the variances of $\hat{\mu}_n^{(2)}$, $\hat{\mu}_n^{(3)}$, and $\hat{\mu}_n^{(4)}$ depend on the unknown parameters (μ, σ^2) .

1. Estimators
2. **Finite Sample Properties**
 - ▷ Bias
 - ▷ Variance
 - ▷ **The Bias-Variance Trade-Off**
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

The Bias-Variance Trade-Off

A popular criterion for evaluating estimators is the mean-squared error:

$$MSE(\hat{\theta}_n) = E\left[\left(\hat{\theta}_n - \theta\right)^2\right]. \quad (8)$$

- ▷ Describes the squared deviations of $\hat{\theta}_n$ from the true parameter.

The next result shows that the MSE is a one-number summary of the bias and variance of an estimator.

Corollary 1

Let $\hat{\theta}_n$ be an estimator for θ . We have

$$MSE(\hat{\theta}_n) = \text{Bias}(\hat{\theta}_n)^2 + \text{Var}(\hat{\theta}_n). \quad (9)$$

Proof.

The Bias-Variance Trade-Off (Contd.)

Example 9

Our analysis suggests that we may prefer $\hat{\mu}_n^{(3)}$ to $\hat{\mu}_n^{(2)}$.

- ▷ Both are unbiased but $\text{Var}(\hat{\mu}_n^{(2)}) > \text{Var}(\hat{\mu}_n^{(3)})$.

But Figure 1 also suggests that we may prefer $\hat{\mu}_n^{(4)}$ to $\hat{\mu}_n^{(2)}$ for small λ .

- ▷ Even though $\text{Bias}(\hat{\mu}_n^{(2)}) < \text{Bias}(\hat{\mu}_n^{(4)})$, we may find the difference in $\text{Var}(\hat{\mu}_n^{(2)})$ and $\text{Var}(\hat{\mu}_n^{(4)})$ sufficiently large to prefer the latter.

Calculations in R show that for the setting of Figure 1, we have:

- ▷ $MSE(\hat{\mu}_n^{(1)}) = 1.00$; $MSE(\hat{\mu}_n^{(2)}) \approx 0.97$;
- ▷ $MSE(\hat{\mu}_n^{(3)}) \approx 0.10$; $MSE(\hat{\mu}_n^{(4)}) \approx 0.09$.

Note: These are results for a *specific parameter values* (μ, σ^2) .
Simulation are not mathematical proofs!

1. Estimators
2. Finite Sample Properties
 - ▷ Bias
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. **Large Sample Properties**
 - ▷ **Consistency**
 - ▷ Asymptotic Distribution
4. On the Interpretation of Estimates

Large Sample Properties

Note that in Examples 5 and 8 depended on unknown parameters (μ, σ^2) .

- ▷ $Bias(\hat{\mu}_n^{(4)})$ depends on μ ;
- ▷ $Var(\hat{\mu}_n^{(2)})$ and $Var(\hat{\mu}_n^{(3)})$ depend on σ^2 ;
- ▷ $Var(\hat{\mu}_n^{(4)})$ depends on (μ, σ^2) .

Without knowledge of the parameters that we want to estimate, we can't rank our estimators in terms of the MSE!

Instead of the (often) impossible question

- ▷ “Which estimator *is* best (or: ‘good enough’)?”

we instead attempt to answer the question

- ▷ “Which estimator *will eventually be* best? (or: ‘good enough’)”

Here, “eventually” considers gathering more and more observations.

Large Sample Properties (Contd.)

It turns out that we can make statements about the *eventual* characteristics of estimators in many settings *without* knowledge of the parameters of interest.

We rely heavily on two notions of convergence of random variables:

- ▷ Convergence in Probability;
- ▷ Convergence in Distribution.

Using these concepts, we study

- ▷ the consistency of an estimator, which checks whether it will eventually be arbitrarily “close” to the true parameter value;
- ▷ the asymptotic distribution of an estimator, which approximates its sampling distribution when n is large.

Convergence in Probability

Recall convergence in the context of sequences of real numbers:

▷ Consider $x, x_1, \dots, x_n \in \mathbb{R}$. We write $x_n \rightarrow x$ if

$$\forall \varepsilon > 0, \exists N_\varepsilon \in \mathbb{N} : |x_n - x| < \varepsilon, \quad \forall n \geq N_\varepsilon.$$

Convergence in probability generalizes this notion of convergence to sequences of random variables.

Definition 2 (Convergence in Probability)

Let X_1, \dots, X_n be a sequence of random variables, and let X be another random variable. We say X_n *converges in probability to* X if

$$\forall \varepsilon > 0, \quad P(|X_n - X| > \varepsilon) \rightarrow 0, \quad \text{as } n \rightarrow \infty. \quad (10)$$

We write $X_n \xrightarrow{P} X$.

In words: If $X_n \xrightarrow{P} X$, then X_n deviates from X by no more than ε with large probability as $n \rightarrow \infty$.

Consistency

We consider convergence in probability to analyze whether an estimator $\hat{\theta}_n$ for θ will eventually be arbitrarily close to the true parameter value.

Definition 3

We say an estimator $\hat{\theta}_n$ for a parameter θ is *consistent* if

$$\hat{\theta}_n \xrightarrow{P} \theta. \quad (11)$$

Consistency is often considered a minimum requirement for an estimator.

- ▷ If the estimator is not arbitrarily close to the true parameter even with infinitely many observations, then there is little hope that it will be reasonably close when the sample size n is finite.
- ▷ *No* inconsistent estimator is considered to be “good enough.”

Note: Equation (11) implicitly considered $n \rightarrow \infty$. Unless otherwise stated, we always consider $n \rightarrow \infty$ in this course.

Example 10

Consider the estimators $\hat{\mu}_n^{(1)}$ and $\hat{\mu}_n^{(2)}$ of Example 5. We have, $\forall \varepsilon > 0$,

$$P\left(|\hat{\mu}_n^{(1)} - \mu| > \varepsilon\right) =$$

$$P\left(|\hat{\mu}_n^{(2)} - \mu| > \varepsilon\right) =$$

Hence, neither $\hat{\mu}_n^{(1)}$ nor $\hat{\mu}_n^{(2)}$ are consistent estimators of μ .

- ▷ Since neither estimator meets the minimum requirement, we won't consider them any further.

Weak Law of Large Numbers

To show consistency of less trivial estimators, we need new technical tools. The most important is the Weak Law of Large Numbers:

Theorem 1 (Weak Law of Large Numbers; WLLN)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} X$ be a random sample. Then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} E[X]. \quad (12)$$

In words: As $n \rightarrow \infty$, the sample average concentrates around its mean.

Example 11

Consider the estimator $\hat{\mu}_n^{(3)}$ of Example 5. By the WLLN,

$$\hat{\mu}_n^{(3)} \xrightarrow{P} \mu,$$

so that $\hat{\mu}_n^{(3)}$ is a consistent estimator of μ .

Weak Law of Large Numbers (Contd.)

To prove the WLLN, we make use of the following intermediate result:

Lemma 1 (Chebyshev's Inequality)

Let X be a random variable. Then,

$$\forall \varepsilon > 0, \quad P(|X| > \varepsilon) \leq \frac{E[X^2]}{\varepsilon^2}. \quad (13)$$

Proof.



Weak Law of Large Numbers (Contd.)

We now return to the proof of the WLLN.

Proof.



Weak Law of Large Numbers (Contd.)

Examples 10 and 11 discussed consistency of the estimators $\hat{\mu}_n^{(1)}$, $\hat{\mu}_n^{(2)}$, and $\hat{\mu}_n^{(3)}$ of Example 5. What about $\hat{\mu}_n^{(4)}$?

Note that

$$\hat{\mu}_n^{(4)} = \frac{1}{n + \lambda} \sum_{i=1}^n X_i = \frac{n}{n + \lambda} \frac{1}{n} \sum_{i=1}^n X_i, \quad (14)$$

so that $\hat{\mu}_n^{(4)}$ is a function of $\frac{1}{n} \sum_{i=1}^n X_i$ and $\frac{n}{n+\lambda}$.

The WLLN provides considers convergence in probability of the sample average. Now, we need tools to:

- ▷ derive convergence in probability of *random vectors*;
- ▷ derive convergence in probability of *functions* of random vectors.

Joint Convergence in Probability

Definition 4

Take $k \in \mathbb{N}$ and let $\tilde{X}_n = (X_{1,n}, \dots, X_{k,n})$, $n \geq 1$, be a sequence of random vectors, and let $\tilde{X} = (X_1, \dots, X_k)$ be another random vector. We say \tilde{X}_n converges in probability to \tilde{X} if

$$\forall \varepsilon > 0, \quad P \left(\sqrt{\sum_{j=1}^k (X_{j,n} - X_j)^2} > \varepsilon \right) \rightarrow 0, \quad \text{as } n \rightarrow \infty. \quad (15)$$

We won't require using Equation (15) directly due to the following result:

Theorem 2

Take $k \in \mathbb{N}$ and let $\tilde{X}_n = (X_{1,n}, \dots, X_{k,n})$, $n \geq 1$, be a sequence of random vectors, and let $\tilde{X} = (X_1, \dots, X_k)$ be another random vector. Then

$$X_{j,n} \xrightarrow{P} X_j, \forall j = 1, \dots, k \quad \Rightarrow \quad \tilde{X}_n \xrightarrow{P} \tilde{X}. \quad (16)$$

Continuous Mapping Theorem

The following theorem delivers a powerful tool for proving convergence of any continuous functions of sample averages.

Theorem 3 (Continuous Mapping Theorem; CMT)

Let $X_n, n \geq 1$, be a sequence of random vectors, and let X be another random vector. If $X_n \xrightarrow{P} X$, then

$$g(X_n) \xrightarrow{P} g(X), \quad (17)$$

for any function g that is continuous at $g(x), \forall x \in \text{supp } X$.

Example 12

Let $A_n \xrightarrow{P} a \in \mathbb{R}$ and $B_n \xrightarrow{P} b \in \mathbb{R}$. Consider $g(a, b) = a/b$. Then

$$g(A_n, B_n) \xrightarrow{P} g(a, b), \quad (18)$$

by the CMT as long as $b \neq 0$.

Continuous Mapping Theorem (Contd.)

Example 13

Consider $\hat{\mu}_n^{(4)}$ from Example 5. We show $\hat{\mu}_n^{(4)} \xrightarrow{P} \mu$ in four steps:

Continuous Mapping Theorem (Contd.)

Example 14

Consider $\hat{\sigma}_n^2$ defined in Example 3. We show $\sqrt{\hat{\sigma}_n^2} \xrightarrow{P} \sigma$ in four steps:

1. Estimators
2. Finite Sample Properties
 - ▷ Bias
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. **Large Sample Properties**
 - ▷ Consistency
 - ▷ **Asymptotic Distribution**
4. On the Interpretation of Estimates

Convergence in Distribution

Examples 11 and 13 showed that both $\hat{\mu}_n^{(3)}$ and $\hat{\mu}_n^{(4)}$ are consistent for θ .

- ▷ But: Consistency does not imply that the choice of estimator is irrelevant even for large n : Could have different variances.

We introduce the concept of convergence in distribution:

- ▷ Allows to assess dispersion of estimators as n grows large.
- ▷ Allows to make approximate probability statements about estimators.

Definition 5 (Convergence in Distribution)

Let $X_n, n \geq 1$, be a sequence of random variables, and let X be another random variable. We say X_n *converges in distribution to* X if

$$P(X_n \leq t) \rightarrow P(X \leq t), \quad \forall t \in \mathbb{R}. \quad (19)$$

We write $X_n \xrightarrow{d} X$.

In words: If $X_n \xrightarrow{d} X$, then the distribution of X_n is approximately equal to the distribution of X for large n .

Central Limit Theorem

The next result is a powerful tool for deriving the asymptotic distribution of sample averages.

Theorem 4 (Central Limit Theorem; CLT)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} X$ be a random sample. Then

$$\frac{\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \right)}{\sigma} \xrightarrow{d} N(0, 1), \quad (20)$$

where $\mu \equiv E[X]$ and $\sigma \equiv sd(X) > 0$.

In words: As n grows large, the distribution of the sample average is approximately normal.

▷ Remarkable because we have *not* assumed that X is normal!

Notation: We could have stated Equation (20) instead as $\frac{\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \right)}{\sigma} \xrightarrow{d} Z$, where $Z \sim N(0, 1)$. As before, we may occasionally use random variables and their distributions interchangeably.

Example 15

Consider $\hat{\mu}_n^{(3)}$ from Example 5. By the CLT, we have

$$\frac{\sqrt{n} \left(\hat{\mu}_n^{(3)} - \mu \right)}{\sigma} \xrightarrow{d} N(0, 1). \quad (21)$$

Hence, for large n , we may approximate the distribution of $\hat{\mu}_n^{(3)}$ with

$$N \left(\mu, \sigma^2/n \right). \quad (22)$$

Note that (22) is of little practical help unless we may substitute parameter estimates for the unknown parameters.

Slutsky's Theorem

Good news: The result of the CLT continues to hold when parameter estimates are substituted for unknown parameter values.

Theorem 5 (Slutsky's Theorem)

Let $A_n, n \geq 1$, and $B_n, n \geq 1$, be sequences of random variables. Let A be another random variable and $b \in \mathbb{R}$. If $A_n \xrightarrow{d} A$ and $B_n \xrightarrow{p} b$, then

$$B_n + A_n \xrightarrow{d} b + A, \quad (23)$$

and

$$B_n A_n \xrightarrow{d} bA. \quad (24)$$

If in addition $b \neq 0$, then also

$$A_n/B_n \xrightarrow{d} A/b. \quad (25)$$

Example 16

Consider $\hat{\sigma}_n^2$ and $\hat{\mu}_n^{(3)}$ from Example 3 and 5. Consider

$$Z_n \equiv \frac{\sqrt{n} \left(\hat{\mu}_n^{(3)} - \mu \right)}{\hat{\sigma}_n} = \frac{\sigma}{\hat{\sigma}_n} \frac{\sqrt{n} \left(\hat{\mu}_n^{(3)} - \mu \right)}{\sigma},$$

so that Slutsky's suggests taking $A_n \equiv \frac{\sqrt{n} \left(\hat{\mu}_n^{(3)} - \mu \right)}{\sigma}$ and $B_n \equiv \frac{\sigma}{\hat{\sigma}_n}$. Then,

Slutsky's Theorem (Contd.)

Example 17

Consider $\hat{\sigma}_n^2$ and $\hat{\mu}_n^{(4)}$ from Example 3 and 5. We want to show that

$$\frac{\sqrt{n} \left(\hat{\mu}_n^{(4)} - \mu \right)}{\hat{\sigma}_n} \xrightarrow{d} N(0, 1).$$

We have

Standard Errors

Informally, Examples 16 and 17 show that the sampling distribution of the estimators can be approximated with $N(\mu, \frac{\hat{\sigma}_n^2}{n})$. For this purpose, practitioners often use so-called *standard errors*.

Definition 6 (Standard Error)

Let $\hat{\theta}_n$ and $\hat{\sigma}_n$ be estimators such that

$$\frac{\sqrt{n} (\hat{\theta}_n - \theta)}{\hat{\sigma}_n} \xrightarrow{d} N(0, 1). \quad (26)$$

The *standard error* of $\hat{\theta}_n$ is defined as

$$se(\hat{\theta}_n) = \frac{\hat{\sigma}_n}{\sqrt{n}}. \quad (27)$$

For large n , we may approximate the sampling distribution of an estimator $\hat{\theta}_n$ for θ with \sqrt{n} -normal asymptotic distribution by $N(\theta, se(\hat{\theta}_n)^2)$.

Confidence Intervals

Researchers often construct asymptotic confidence intervals to succinctly characterize the approximate sampling distribution:

Theorem 6

Let $\hat{\theta}_n$ be an estimator for θ such that (26) holds. For $\alpha \in (0, 1)$, consider

$$C_n = \left[\hat{\theta}_n - z_{1-\frac{\alpha}{2}} \text{se}(\hat{\theta}_n), \quad \hat{\theta}_n + z_{1-\frac{\alpha}{2}} \text{se}(\hat{\theta}_n) \right], \quad (28)$$

$$C_n^+ = \left[\hat{\theta}_n - z_{1-\alpha} \text{se}(\hat{\theta}_n), \quad \infty \right), \quad (29)$$

$$C_n^- = \left(-\infty, \quad \hat{\theta}_n + z_{1-\alpha} \text{se}(\hat{\theta}_n) \right], \quad (30)$$

where $z_{1-a} \equiv \Phi^{-1}(1-a)$ is the $1-a$ quantile of a standard normal.

C_n, C_n^+ and C_n^- are asymptotically valid $1-\alpha$ confidence intervals. I.e.,

$$P(\theta \in \tilde{C}_n) \rightarrow 1 - \alpha, \quad (31)$$

for $\tilde{C}_n = C_n, C_n^+, C_n^-$.

Confidence Intervals (Contd.)

Proof.

We prove the theorem only for the symmetric confidence interval C_n .
Proofs for C_n^+ and C_n^- are left as a self-study exercise.



Bivariate Central Limit Theorem

Slutsky's Theorem considered the joint convergence of sequences of random variables when one of the sequences converges to a constant.

- ▷ Need tools to understand joint convergence when *both* sequences converge to a random variable. Fortunately, we have the next result:

Theorem 7 (Bivariate Central Limit Theorem)

Let $\tilde{X}_1, \dots, \tilde{X}_n \stackrel{iid}{\sim} Y$ be a sample of bivariate random vectors where $\tilde{X}_i = (X_{1,i}, X_{2,i})$ and $\tilde{X} = (X_1, X_2)$. Then

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \tilde{X}_i - \mu \right) \xrightarrow{d} N(0, \Sigma), \quad (32)$$

where $\mu \equiv E[\tilde{X}]$ and

$$\Sigma \equiv \text{Var}(\tilde{X}) = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{bmatrix}. \quad (33)$$

Example 18

Consider a sample $(Y_1, X_1), \dots, (Y_n, X_n) \stackrel{iid}{\sim} (Y, X)$ where $X \sim \text{Bernoulli}(p)$ with unknown $p \in (0, 1)$. Suppose we are interested in the joint distribution of the estimators

$$E_n[YX] = \frac{1}{n} \sum_{i=1}^n Y_i X_i, \quad \text{and} \quad E_n[Y(1-X)] = \frac{1}{n} \sum_{i=1}^n Y_i (1 - X_i). \quad (34)$$

By the (bivariate) CLT, we have

Bivariate Slutsky's Theorem

As was the case with the univariate CLT, it's bivariate analogue is particularly useful when combined with a Slutsky-type result:

Theorem 8 (Bivariate Slutsky's Theorem)

Let $A_n, n \geq 1$, and $B_n, n \geq 1$, be sequences of bivariate random vectors variables. Let A be another bivariate random vector and $b \in \mathbb{R}^2$. If $A_n \xrightarrow{d} A$ and $B_n \xrightarrow{p} b$, then

$$A_n + B_n \xrightarrow{d} A + b, \quad (35)$$

and

$$B_n^\top A_n \xrightarrow{d} b^\top A. \quad (36)$$

Bivariate Slutsky's Theorem (Contd.)

Example 19

Let $A_n, n \geq 1$ and $B_n, n \geq 1$ be sequences of bivariate random vectors such that $A_n \xrightarrow{d} N(0, \Sigma)$ and $B_n \xrightarrow{P} b \in \mathbb{R}^2$. By Slutsky's Theorem,

$$B_n^\top A_n \xrightarrow{d} b^\top N(0, \Sigma) \stackrel{d}{=} N(0, b^\top \Sigma b),$$

where the last equation follows from Lemma 4c of Lecture 2A.

Suppose now that $Z_n, n \geq 1$, such that $Z_n \xrightarrow{d} N(0, I_2)$, and $\hat{\Sigma}_n, n \geq 1$ is a sequence of estimators such that $\hat{\Sigma}_n^{-1}$ exists and $\hat{\Sigma}_n \xrightarrow{P} \Sigma$. By the CMT,

whenever Σ^{-1} exists. Hence, by Slutsky's Theorem,

Example 20

Consider the setting of Example 18 and construct the estimator

$$E_n[YX] - E_n[Y(1 - X)] = \begin{bmatrix} 1 \\ -1 \end{bmatrix}^\top \begin{bmatrix} E_n[YX] \\ E_n[Y(1 - X)] \end{bmatrix}. \quad (37)$$

Hence, it follows from Example 18 and Slutsky's Theorem that

1. Estimators
2. Finite Sample Properties
 - ▷ Bias
 - ▷ Variance
 - ▷ The Bias-Variance Trade-Off
3. Large Sample Properties
 - ▷ Consistency
 - ▷ Asymptotic Distribution
4. **On the Interpretation of Estimates**

On the Interpretation of Estimates

Thus far, we have exclusively discussed *estimators* $\hat{\theta}_n$ for a parameter θ .

- ▷ $\hat{\theta}_n$ is a function of the sample $X_1, \dots, X_n \sim X$ and thus random.
- ▷ How does real-world data come in?

Data is a realization of our sample X_1, \dots, X_n .

- ▷ The data we have collected is the collection of numbers: x_1, \dots, x_n .

An *estimate* is a realization of our estimator $\hat{\theta}_n$:

- ▷ The estimator $\hat{\theta}_n(X_1, \dots, X_n)$ is a random variable;
- ▷ The estimate $\hat{\theta}_n(x_1, \dots, x_n)$ is a number.

This distinction between estimators and estimates can lead to confusion.

- ▷ We can make probabilistic statements about $\hat{\theta}_n(X_1, \dots, X_n)$.
- ▷ We cannot make probabilistic statements about $\hat{\theta}_n(x_1, \dots, x_n)$.

Notation: To make matters worse, $\hat{\theta}_n$ often denotes both the estimator (random) and the estimate (fixed), so that you have to figure it out yourself from context!

On the Interpretation of Estimates (Contd.)

The confusion between estimators (random) and estimates (fixed) is particularly severe in the context of confidence intervals.

Recall that an asymptotic $1 - \alpha$ confidence interval is such that

$$P(\theta \in C_n) \rightarrow 1 - \alpha.$$

Let c_n denote a realization of C_n (i.e., what you computed using data).

- ▷ It is correct to say C_n covers θ w.p. (tending to) $1 - \alpha$.
- ▷ It is *incorrect* to say c_n covers θ w.p. (tending to) $1 - \alpha$.
- ▷ $P(\theta \in c_n) = \mathbb{1}\{\theta \in c_n\} \in \{0, 1\}$. This is a comparison of numbers!

On the Interpretation of Estimates (Contd.)

Statistics courses often introduce the idea of repeated experiments to interpret confidence intervals:

- ▷ “If I were to repeat the same experiment again and again, each time computing a $1 - \alpha$ confidence interval, then the confidence intervals would cover the true parameter $100(1 - \alpha)\%$ of the time.”

This interpretation is correct but requires some mental gymnastics.

- ▷ The same experiment is seldom repeated many times.
- ▷ Only hypothetically reassuring.

A more useful interpretation is the idea of many unrelated experiments:

- ▷ In your career, you are going to calculate many $1 - \alpha$ confidence intervals for unrelated parameters. Of these confidence intervals, $100(1 - \alpha)\%$ cover the corresponding true parameter.

What does this tell you about whether a specific θ is in a computed c_n ?

- ▷ Nothing! You are correct $100(1 - \alpha)\%$ of the time, but you'll never know when. (If you're discontent: Check out *Bayesian* statistics.)

Example 21

Consider $\hat{\mu}_n^{(3)}$ from Example 5. Suppose that we collected data and that

- ▷ $\hat{\mu}_n^{(3)} = 10$;
- ▷ $se(\hat{\mu}_n^{(3)}) = 3$.

Then, an asymptotic $1 - \alpha$ confidence interval given by

$$c_n = \tag{38}$$

Here c_n denotes a realization of the confidence interval C_n :

- ▷ We've collected data (a realization of our sample);
- ▷ Computed the estimator $\hat{\mu}_n^{(3)}$ and its standard error $se(\hat{\mu}_n^{(3)})$;
- ▷ Calculated a $1 - \alpha$ confidence interval c_n .

What is $P(\theta \in c_n)$?

On the Interpretation of Estimates (Contd.)

Example 22

The GRE is a standardized test required for admission to many graduate programs in the US. Test-takers receive three scores (Verbal, Writing, and Quantitative). Below is a screenshot from their documentation.

Table 5A: Reliability Estimates and Standard Errors of Measurement (SEM)^a
for Individual Scores and Score Differences for the GRE[®] General Test

Score	Reliability Estimate	SEM of Individual Scores	SEM of Score Differences
Verbal Reasoning	0.93	2.4	3.4
Quantitative Reasoning	0.95	2.2	3.1
Analytical Writing	0.87	0.30	0.43

(Don't get confused by the different terminology: SEMs are essentially the same as the standard errors discussed earlier.)

Suppose that a student received a quantitative score of 153 (out of 170). Then a 95% confidence interval for her “true score” is given by

$$c_n = \quad (39)$$

Example 22 (Contd.)

How difficult is it to get the interpretation of asymptotic confidence intervals correct?

- ▷ It turns out: Difficult enough for the GRE...

In the “GRE Guide to the Use of Scores” ([link](#)), they write:

- ▷ “consider a test taker who obtained a GRE Quantitative test score of 153. [...] we can be 95% confident that the test taker’s true score would be between 149 and 157.”

This is precisely the incorrect interpretation of confidence intervals!

What makes this more amusing is that the title of the corresponding section is “Interpret GRE Scores Carefully [...]”.

(We all need a little bit of joy in our life, but typically it’s appreciated when you don’t make too much fun of others’ statistics skills.)

Summary

This concludes the Part A of our statistics review.

- ▷ Introduced the sample analogue principle to develop estimators;
- ▷ Discussed finite sample properties of estimators, in particular, their bias, variance, and MSE;
- ▷ Generalized the concept of convergence to random variables via convergence in probability and convergence in distribution;
- ▷ Studied large sample properties of estimators, in particular, their consistency and asymptotic distribution.

A key insight was that under fairly general conditions, approximate probabilistic statements about estimators can be made using their asymptotic distribution.

- ▷ In Part B, we discuss how estimators and their (approximate) sampling distributions can be leveraged to assess whether the true parameter θ takes a particular value, say, θ_0 .
- ▷ This is known as *hypothesis testing*.

References

Wasserman, L. (2003). *All of statistics*. Springer.