

Question 2: Confidence Intervals

Setup: $\{X_i\}_{i=1}^N$ iid

$E[X_i] = M$ unknown

$\text{Var}(X_i) = \sigma^2$ known finite

Principle: construct a test statistic,
find its distribution, and invert.

find the values such that a test doesn't reject

CI for M : example we can write $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$
sample mean.

1. Find CI $[a, b]$ such that

$$P(M \in [a, b]) \rightarrow 1 - \alpha \text{ as } N \rightarrow \infty$$

Ans: use $t_n = \frac{\bar{X} - M}{\sigma/\sqrt{n}}$ (t-statistic)

$$E t_n = \frac{1}{\sigma\sqrt{n}} \left(\underbrace{\sum_i E[X_i] - M}_n \right) = 0$$

$$\text{Var}(t_n) = \left(\frac{1}{\sigma\sqrt{n}} \right)^2 \frac{1}{n} \sum_i \text{Var}(X_i) = 1$$

Apply CLT: $t_n \xrightarrow{d} Z \sim N(0, 1)$

CDF of $t_n \rightarrow$ CDF of Z pointwise

Condition acceptance region $(-q_{1-\alpha/2}, q_{\alpha/2})$

where q_α is the α^{th} quantile of Z

$$\Pr(Z \in [-q_{1-\alpha}, q_{\alpha}]) = 1 - \alpha$$

So $\Pr(t_n \in [-q_{1-\alpha}, q_{\alpha}]) \rightarrow 1 - \alpha$ as $N \rightarrow \infty$

From central limit theorem

$$\begin{aligned} & \Pr\left(\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \in [-q_{1-\alpha}, q_{\alpha}]\right) \\ &= \Pr\left(-\frac{\sigma}{\sqrt{n}} q_{1-\alpha} \leq \bar{x} - \mu \leq \frac{\sigma}{\sqrt{n}} q_{\alpha}\right) \\ &= \Pr\left(\bar{x} + q_{\alpha} \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \bar{x} + q_{\alpha} \frac{\sigma}{\sqrt{n}}\right) \end{aligned}$$

$$\begin{array}{ll} \text{if } \alpha = 0.05, & q_{\alpha} = 1.96 \\ \text{if } \alpha = 0.01, & q_{1-\alpha} = 2.33 \end{array}$$

Known vs unknown or does it change this result, because of Slutsky's theorem

2. Suppose goal is $[a, b]$ such that

$$\Pr(M \in [a, b]) \geq 1 - \alpha$$

This is called a "Probably Approximately Correct" (PAC) statistic. Common in ML and comp sci.

→ Cetin Valiant : popular book
Method: Chebychev's inequality

$$\Pr\left(\left|Z - \mathbb{E}(Z)\right| \geq k\sqrt{\text{Var}(Z)}\right) \leq \frac{1}{k^2}$$

Approach: use $t_n = \frac{\bar{x} - M}{\sigma/\sqrt{n}}$: $\mathbb{E}[t_n] = 0$
 $\text{Var}(t_n) = 1$

Apply Chebychev for t_n

$$\Pr\left(\left|t_n\right| \geq k\right) \leq \frac{1}{k^2}$$

so from 1-d interval, $\alpha = \frac{1}{k^2}$

$$\text{Set } K = \alpha^{-\frac{1}{2}}$$

so test: accept if $|t_n| \leq \alpha^{-\frac{1}{2}}$, reject otherwise.

Invert 1-qf interval

$$[\bar{a}, \bar{b}] = \left[\bar{x} - \alpha^{-\frac{1}{2}} \frac{\sigma}{\sqrt{n}}, \bar{x} + \alpha^{-\frac{1}{2}} \frac{\sigma}{\sqrt{n}}\right]$$

$$\text{hence } \Pr(M \in [\bar{a}, \bar{b}]) \geq 1 - \alpha$$

? Why don't people do this?

Reason is, $[\bar{a}, \bar{b}]$ in (2) is much larger than $[a, b]$ in (1).

Why care? Power, if we want to decide based on
 Want to balance Type I & Type II errors
 If interval is large, size \downarrow , power \downarrow also,
 \Rightarrow will make fewer type I, more type II errors.

Can see this two ways: Normal distribution converges to
 has quantiles which grow logarithmically
 with λ , while Chibyshev bound grows
 quadratically, i.e., especially for small λ
 Chibyshev bound is much larger. Ratio of size
 goes to ∞ as $\lambda \rightarrow 0$

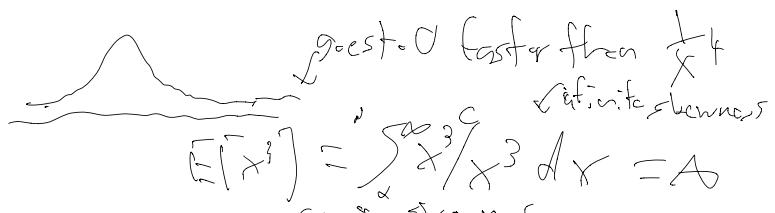
Ex: $\lambda = 0.01$, for normal, this is 2.6
 for Chibyshev, $\lambda^{-\frac{1}{2}} = 10$, so interval is
 ≈ 4 times as wide.

This would be reasonably okay if bound were sharp,
 and actual probability was close to 1 - λ
 but if inequality is slack
 $\Pr(M \in [a, b]) \gg (\lambda)$,
 then might suffer from redundancy just like before.

Part 4) Get a closer bound with more assumptions.
 N, finite sample rates. If just finitely varying with CLT
 CDF converges, but no idea how fast.

And assumption: $E|x_i - \mu|^3 < \infty$ finite skewness,

fail. + pdf



S. this limitation
on probability of
large outliers

$$\begin{aligned} E[x^3] &= \int_0^\infty x^3 / x^{4+\epsilon} dx \\ &= \int_0^\infty x^{1-\epsilon} dx < \infty \end{aligned}$$

(can occur in financial data, your average is mostly determined by Black Monday, 1987 instead 500
for decades after, 1 day determining average returns)

Methd for f : used f convgaes in distribution

$$Ef(T_n) \rightarrow Ef(S_n) \quad \text{convgaes in f}$$

$$\begin{aligned} S_n &= \frac{1}{n} \sum_{i=1}^n \frac{x_i - \mu}{\sigma} & T_n &= \sum_{i=1}^n z_i \\ \mathbb{E} \frac{x_i - \mu}{\sigma} &= 0 & z_i &\stackrel{iid}{\sim} N\left(0, \frac{\sigma^2}{n}\right) \\ \text{var}\left(\frac{x_i - \mu}{\sigma}\right) &= \frac{\sigma^2}{n} \end{aligned}$$

What needs: for (\mathbb{E}, f) : $\mathbb{E}\{I\{x \in [a, b]\}\}$
if we consider

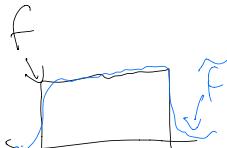
$$|\mathbb{E}f(T_n) - \mathbb{E}f(S_n)| \leq \text{bound}$$

then can use f. calcult. (a, b) with desired probability,

Trick wic Taylorapprox. the input of CLT.

but f is not C^3 , so we have,

$\forall \varepsilon > 0, \exists \tilde{f}$ such that $\tilde{f} \in C^3$ and



$$|\mathbb{E}[f(X)] - \mathbb{E}[\tilde{f}(X)]| < \varepsilon$$

$$|\mathbb{E}[f(S_n)] - f(\bar{T}_n)|$$

$$\leq |\mathbb{E}[f(S_n) - \tilde{f}(S_n)]| < \varepsilon$$

$$+ |\mathbb{E}[\tilde{f}(S_n) - \tilde{f}(\bar{T}_n)]|$$

$$+ |\mathbb{E}[\tilde{f}(\bar{T}_n)] - \mathbb{E}[\tilde{f}(\bar{T})]|$$

$$< \varepsilon$$

$$\leq 2\varepsilon + |\mathbb{E}[\tilde{f}(S_n) - \tilde{f}(\bar{T}_n)]|$$

$$\underbrace{\text{With } n \rightarrow 0}_{\text{it's zero}} \leq \frac{1}{\sqrt{n}} \mathbb{E}[(\bar{T}_n - M)^3]$$

$$\leq \frac{16}{\sqrt{n}}$$

«final»

$$= \frac{16}{\sqrt{n}}$$

»

So apply Triangular Inequality

$$\Pr(S_n \in [a, b]) \geq \Pr(\bar{T}_n \in [\alpha, \beta]) - \frac{16}{\sqrt{N}} \quad \text{making } \geq \text{d}$$

But we know distribution of \bar{T}_n exactly. Find values a and b such that $d = b - a$.

$$\bar{T}_n \sim N(1, \frac{1}{N}) \quad \text{so,} \quad a = \sigma \Phi^{-1}\left(\frac{1}{2}(\alpha - \frac{16}{\sqrt{N}})\right)$$
$$b = \sigma \Phi^{-1}\left(\frac{1}{2}(1 - \frac{1}{2}(\alpha - \frac{16}{\sqrt{N}}))\right)$$

$$\text{Then } \Pr(S_n \in [a, b]) \geq \Pr(\bar{T}_n \in [\alpha, \beta]) - \frac{16}{\sqrt{N}} - \frac{16}{\sqrt{N}} \quad \text{= d}$$

So, (\bar{T}_n) by inverting is

$$\left[\bar{T}_n - \frac{\sigma}{\sqrt{N}} \Phi^{-1}\left(\frac{1}{2}(\alpha - \frac{16}{\sqrt{N}})\right), \right.$$

$$\left. \bar{T}_n - \frac{\sigma}{\sqrt{N}} \Phi^{-1}\left(\frac{1}{2}(1 - \frac{1}{2}(\alpha - \frac{16}{\sqrt{N}}))\right) \right]$$

use slightly different quantiles of normal distribution
whose quantile depends on N and gets closer to α
as $N \rightarrow \infty$

1 - contains same as CLT, rate of width +
CLT interval approx. 1 as N $\rightarrow \infty$

This kind of result follows from central CLT
limit distributions. (PAC)

- Rostamizadeh, M. Er. Tukwalkar: "No guarantees
Archimedes for OLS"

Recently, Nasheen done, kanchibhotla et al
applied this kind of reasoning to OLS in linear
without assuming normality or exact linear fit.
OLS converges to $(\sum x_i x_i^T)^{-1} \sum x_i y_i$

2. OLS with multipl. Regressors

$$Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon \quad X_1 \in \mathbb{R}^{N \times K_1}, \\ X_2 \in \mathbb{R}^{N \times K_2}$$

2.1. What if we don't see X_2 ?
What can we say about $\hat{\beta}_1$?

Sufficient conditions for unbiased estimate,

Try estimator as with just X_1
 model: $\hat{Y} = \hat{X}_1 \hat{\beta}_1 + \hat{U}$

const + \hat{Y} try obs: $(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{Y} = \hat{\beta}_1$
 Writings

$$\begin{aligned}\hat{\beta}_1 &= \underbrace{(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{X}_1 \beta_1}_{I} + (\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{U} \\ &= \beta_1 + (\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{X}_1 \beta_1 + (\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{U} \\ S \sim E[\hat{\beta}_1 - \beta_1] &= E[(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{X}_1 \beta_1 + (\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{U}] \\ &= 0\end{aligned}$$

Visual condition, $E[\hat{U}|X_1, X_2] = 0$ is needed for unbiasedness.

$$\begin{aligned}\text{implies } E[(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{U}] &= E[(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \underbrace{E[\hat{U}|X_1, X_2]}_0] \\ &= 0\end{aligned}$$

have condition

$$E[(\hat{X}_1 \hat{X}_1)^{-1} \hat{X}_1 \hat{X}_2 \beta_2] = 0$$

When is this true?

Case 1: $E[X_2 | X_1] = 0$

Then conditioning on X_1 , get $E(X_2 | X_1) - \{X_2 - E[X_2 | X_1]\}$

Case 2, first example of orthogonality
 $X_2 - \{X_2 - E[X_2 | X_1]\} \perp \text{orthogonal}$

Geometrically, if there are no right angles, \checkmark if X_2 and X_1 are uncorrelated.

In other words if X_1 and X_2 are uncorrelated,

Also fine if $B_2 = 0$

Related to omitted variables bias, if we want b_1 ,

Worry about X_2 if it's correlated with X_1
and $B_2 \neq 0$

2, New estimators, show it is unbiased,

a) Regress y on X_1 , get residuals $\hat{\epsilon}_1$

b) Regress each column of X_2 on X_1 , get residuals \hat{X}_2

c) Regress $\hat{\epsilon}_1$ on \hat{X}_2 + $\text{st } B_2$

Result: a) Residuals $\hat{e}_i = \gamma - \mathbf{x}_i^T (\mathbf{x}_i^T \mathbf{x}_i)^{-1} \mathbf{x}_i^T \gamma$

$$= \underbrace{(\mathbf{I} - \mathbf{x}_i^T (\mathbf{x}_i^T \mathbf{x}_i)^{-1} \mathbf{x}_i^T)}_{M_i} \gamma$$

M_i is orthogonal projection onto complement span of \mathbf{x}_i

b) Parameter $\hat{\gamma}_k$:
Residual \hat{x}_i^k , each column \hat{x}_i^k on \mathbf{x}_i

$$\hat{x}_i^k = (\mathbf{I} - \mathbf{x}_i^T (\mathbf{x}_i^T \mathbf{x}_i)^{-1} \mathbf{x}_i^T) \mathbf{x}_i^k$$

$$= M_i \mathbf{x}_i^k$$

(in assembling matrix $[\hat{x}_1^1 \ \hat{x}_2^1 \ \dots \ \hat{x}_n^1]$)

$$= \hat{x}_i = M_i \mathbf{x}_i$$

\mathbf{x}_i (projection onto span($\mathbf{x}_1, \mathbf{x}_2$)) c) Regress \hat{e}_i on \hat{x}_i

Formula $(\hat{x}_i^T \hat{x}_i)^{-1} \hat{x}_i^T \hat{e}_i = \hat{B}_i$

$$= (\hat{x}_i^T M_i M_i \hat{x}_i)^{-1} \hat{x}_i^T M_i M_i \gamma$$

$$= (\hat{\beta}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2 Y$$

by idempotence,

Therefore we formulate $Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$

$$\hat{\beta}_2 = (\hat{\chi}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2 \underbrace{(X_1 \beta_1 + X_2 \beta_2 + \epsilon)}_0$$

$M_2 \hat{\chi}_2 = 0$ by construction

$$= \underbrace{(\hat{\chi}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2}_{E} \underbrace{\hat{\chi}_2' \beta_2}$$

$$+ \underbrace{(\hat{\chi}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2 \epsilon}$$

$$S.E(\hat{\beta}_2 - \beta_2 | \hat{\chi}_1, \hat{\chi}_2)$$

$$= E[(\hat{\chi}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2 \underbrace{S.E(\hat{\chi}_1, \hat{\chi}_2)}_{0 \text{ unstructured}}]$$

$$S.E(\hat{\beta}_2 - \beta_2 | \hat{\chi}_1, \hat{\chi}_2)$$

regression

$$\text{formula: } \hat{\beta}_2 = (\hat{\chi}_2' M_2 \hat{\chi}_2)^{-1} \hat{\chi}_2' M_2 Y$$

It's called the partitioned regression formula.

It is equal to $\hat{\beta}_j$ component-facs.

This estimator is consistently OLS.

This is known as Frisch-Waugh-Lovell theorem.

3. Linear Probability Model

$$\{Y_i, X_i\}_{i=1}^N \quad Y_i \in \{0, 1\}$$

$$\text{Mod. } P(Y_i = 1 | X_i) = X_i^\top \beta$$

(Note: if X continuous, have the constraint that
third-cointg. above 1 or below 0)

Estimating equation: $\hat{Y}_i = X_i^\top \beta + \varepsilon_i$

$$\text{Note } P(Y_i = 1 | X_i)$$

$$= E[Y_i | X_i] = 1 \cdot P(Y_i = 1 | X_i) +$$

$$\underbrace{0 \cdot P(Y_i = 0 | X_i)}_{0}$$

So can apply results on conditional expectation.

If $E[Y_i | X_i] = X_i^\top \beta$, then OLS is unbiased.

$$\varepsilon_i = y_i - \hat{x}_i \beta$$

$$E[\varepsilon_i | \hat{x}_i] = E[y_i(\hat{x}_i) - \hat{x}_i \beta]$$

$$= \hat{x}_i \beta - \hat{x}_i \beta$$

$$= 0$$

Why does N.F. B(VE)?

Need for Gauss Markov homoskedasticity

$$\text{Var}(\varepsilon_i | \hat{x}_i)$$

$$= E[\varepsilon_i^2 | \hat{x}_i] = \underbrace{E[(\varepsilon | \hat{x})^2]}_0$$

$$= E((y_i - \hat{x}_i \beta)^2 | \hat{x}) = E[y^2 | \hat{x}] + E[\hat{x}^2 \beta^2 | \hat{x}]$$

$$- 2E[y \hat{x} \beta | \hat{x}]$$

N.F. $\hat{y} = y$ since $y \sim \sigma \mathcal{N}$

so this

$$E[y^2 | \hat{x}] + E[E[\hat{x}^2 | \hat{x}]^2 | \hat{x}] - 2E[E[\hat{x}^2 | \hat{x}] | \hat{x}]$$

$$= E[y^2 | \hat{x}] - E[y | \hat{x}]^2$$

$$= \hat{x}^2 \beta^2 - (\hat{x} \beta)^2 = \hat{x} \beta (1 - \hat{x} \beta)$$

is a function of \hat{x} ,

$$= p_r(y_i = 1 | \hat{x}) (1 - p_r(y_i = 1 | \hat{x}))$$

(binomial variance form) ;
 Consider instead GLS:

$$\hat{\beta}^{GLS} = \underbrace{(\mathbf{x}' \text{Var}(\epsilon|\mathbf{x})^{-1} \mathbf{x})^{-1}}_{\Sigma} \mathbf{x}' (\text{Var}(\epsilon|\mathbf{x})^{-1}) \mathbf{y}$$

$$\sigma_{\epsilon_i}^2 \Sigma_{ii} = (\mathbf{x}' \hat{\beta})(1 - \mathbf{x}' \hat{\beta})$$

$$\epsilon_{ij} = \sigma \quad \text{for } i \neq j \quad \text{and} \quad$$

But this depends on true β ,

few steps down, get $\hat{\beta}^{OLS}$ from $y = x\beta + \epsilon$
 Get $\hat{\sigma}_{\epsilon}^2 = \mathbf{x}' \hat{\beta} (1 - \mathbf{x}' \hat{\beta}) \rightarrow \hat{\sigma}_{\epsilon}^2$ by consistency
 $\Sigma = \text{diag}(\hat{\sigma}_{\epsilon}^2)$

$$\text{thus } \hat{\beta}^{FGLS} = (\mathbf{x}' \Sigma^{-1} \mathbf{x})^{-1} \mathbf{x}' \Sigma^{-1} \mathbf{y}$$

asymptotically efficient by efficiency result for GLS.

Here assumed correct specification of error variance,
 which for LPM is probably wrong
 [linear probit, logit model]