

Convergence concept.

$\{X_n\}_{n=1}^{\infty}$ probability space (Ω, \mathcal{F}, P) .

$X_n: \Omega \rightarrow \mathbb{R}$. $X_n(\omega)$

Def: $X_n \rightarrow X$ a.s.

$$P(\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\}) = 1$$

$$\text{or: } P(\lim_{n \rightarrow \infty} X_n = X) = 1$$

most generally: $P(\lim_{n \rightarrow \infty} d(X_n, X) = 0) = 1$ for random element X_n .

Def: convergence in mean square.

$$r \in (0, +\infty), \lim_{n \rightarrow \infty} E[|X_n - X|^r] = 0$$

for $r=2$: proposition: $\lim_{n \rightarrow \infty} E[X_n] = c$

$X_n \rightarrow c$ in mean square iff $\lim_{n \rightarrow \infty} \text{Var}(X_n) = 0$.

Theorem: Portmanteau Lemma:

$$\textcircled{1} X_n \xrightarrow{d} X \Leftrightarrow \textcircled{2} \lim_{n \rightarrow \infty} P(X_n \in A) = P(X \in A).$$

Theorem: Algebra.

$$\textcircled{1} X_n \xrightarrow{P} X, Y_n \xrightarrow{P} Y \Rightarrow X_n + Y_n \xrightarrow{P} X + Y$$

$$\textcircled{2} X_n \rightarrow X \text{ a.s.}, Y_n \rightarrow Y \text{ a.s.} \Rightarrow X_n + Y_n \rightarrow X + Y \text{ a.s.}$$

\textcircled{3} $X_n \rightarrow X$ in r-mean, $Y_n \rightarrow Y$ in r-mean \Rightarrow same.

$$\textcircled{4} X_n \xrightarrow{d} X, Y_n \xrightarrow{d} Y, X_n \perp Y_n \text{ f.n., } X \perp Y.$$

$$X_n + Y_n \xrightarrow{d} X + Y$$

$$\textcircled{5} X_n \xrightarrow{d} X, Y_n \xrightarrow{d} \alpha, \alpha \in \mathbb{R} \Rightarrow X_n + Y_n \xrightarrow{d} X + \alpha$$

Def: $X_n \xrightarrow{P} X$:

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$$

$$\text{or } \lim_{n \rightarrow \infty} P(|X_n - X| \leq \varepsilon) = 1$$

most generally: $\lim_{n \rightarrow \infty} P(d(X_n, X) \geq \varepsilon) = 0$

Def: $X_n \xrightarrow{d} X$:

$$\text{if: } F_{X_n}(x) \rightarrow F_x(x) \quad \forall x \in C(F_x).$$

limit distribution function

Relationships:

a.s. \Rightarrow in prob. \Rightarrow in dist.

\uparrow
in r-mean

$X_n \xrightarrow{d} \alpha \Rightarrow X_n \xrightarrow{P} \alpha$

\textcircled{6} $\cdot X_n \xrightarrow{P} X, Y_n \xrightarrow{P} Y$: then,

$$X_n - Y_n \xrightarrow{P} X - Y \quad X_n Y_n \xrightarrow{P} X Y$$

if $P(X_n \neq 0) = 1$ and $P(Y_n \neq 0) = 1$, then $\frac{X_n}{Y_n} \xrightarrow{P} \frac{X}{Y}$

for r.vectors: \textcircled{7} $X_n \xrightarrow{P} X, Y_n \xrightarrow{P} Y \Rightarrow (X_n, Y_n) \xrightarrow{P} (X, Y)$

\textcircled{8} $X_n \rightarrow X \text{ a.s.}, Y_n \rightarrow Y \text{ a.s.} \Rightarrow (X_n, Y_n) \rightarrow (X, Y) \text{ a.s.}$

\textcircled{9} $X_n \xrightarrow{d} X, Y_n \rightarrow c, c \in \mathbb{R}, \Rightarrow (X_n, Y_n) \xrightarrow{d} (X, c)$.

The continuity theorem and continuous mapping theorem.

Theorem: Lévy's continuous theorem:

$$\{X_n\}: r.v. \phi: X_n \rightarrow \phi X \text{ pointwise on all } \mathbb{R} \Leftrightarrow X_n \xrightarrow{d} X$$

unformular.

Stochastic Order Notation:

Def: $X_n = O_p(Y_n)$:

$$\{X_n\} \text{ r.v. } \{Y_n\} \in \mathbb{R}, Y_n > 0$$

$\forall \varepsilon > 0, \exists M_\varepsilon \text{ and } N_\varepsilon, \text{ for all } n > N_\varepsilon$,

$$P\left(\left|\frac{X_n}{Y_n}\right| \geq M_\varepsilon\right) = P(|X_n| \geq Y_n M_\varepsilon) \leq \varepsilon$$

$\Leftrightarrow P\left(\left|\frac{X_n}{Y_n}\right| \leq M_\varepsilon\right) \geq 1 - \varepsilon; \frac{X_n}{Y_n}$ is bounded in probability

Def: $X_n = o_p(Y_n)$:

$$\{X_n\}: r.v. \{Y_n\} \in \mathbb{R}, Y_n > 0$$

$\left|\frac{X_n}{Y_n}\right| \xrightarrow{P} 0 \text{ as } n \rightarrow \infty$

$X_n \xrightarrow{P} 0 \text{ at the rate } Y_n$.

Theorem: $\{Y_n\}$ r.v.

$$\textcircled{1} o_p(u) + o_p(v) = o_p(u+v)$$

$$\textcircled{2} o_p(u) + o_p(cu) = o_p(u)$$

$$\textcircled{3} o_p(u) o_p(v) = o_p(uv)$$

$$\textcircled{4} (1 + o_p(1))^{-1} = o_p(1)$$

$$\textcircled{5} o_p(cu) = cu o_p(u)$$

$$\textcircled{6} o_p(uv) = u o_p(v)$$

$$\textcircled{7} o_p(o_p(u)) = o_p(u)$$

Theorem: $\hat{\theta}_n$: r.v. b. scalar.

suppose $\hat{\theta}_n(\hat{\theta}_n - \theta) \xrightarrow{d} Z$ for some $\hat{z} > 0$, then $\hat{\theta}_n \xrightarrow{P} \theta$.

proof: $\hat{\theta}_n(\hat{\theta}_n - \theta) \xrightarrow{d} Z \Rightarrow \hat{\theta}_n(\hat{\theta}_n - \theta) = o_p(1)$, thus, $\hat{\theta}_n - \theta = \frac{1}{\hat{\theta}_n} \hat{\theta}_n(\hat{\theta}_n - \theta) = \frac{1}{\hat{\theta}_n} o_p(1) = o_p(1) \xrightarrow{P} 0 \Rightarrow \hat{\theta}_n \xrightarrow{P} \theta$.

LLN:

Theorem: Bernoulli's LLN:

Let $X_1 \dots X_n$ iid Bernoulli r.v. with p . $\bar{X}_n \xrightarrow{P} p$

Theorem: Weak LLN:

$\{X_i\}_{i=1}^{\infty}$ r.v. suppose: \textcircled{1} independent $\{X_i\}$ \textcircled{2} $E[X_i] = \mu < \infty, \forall i$ \textcircled{3} $\exists M > 0, \text{Var}(X_i) \leq M, \forall i$

then: $\bar{X}_n \xrightarrow{P} \mu$

pf: $E[\bar{X}_n] = \mu, \text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \leq \frac{M}{n}$. $P(|\bar{X}_n - E[\bar{X}_n]| \geq \varepsilon) \leq \frac{\text{Var}(\bar{X}_n)}{\varepsilon^2} \leq \frac{1}{n} \frac{M}{\varepsilon^2} \rightarrow 0$ as $n \rightarrow \infty$.

Theorem: Kolmogorov's 2nd strong LLN: most important, strongest.

$\{X_i\}$: iid r.v., then $\bar{X}_n \rightarrow \mu$ a.s. as $n \rightarrow \infty$ iff $E[X_i] = \mu < \infty$ exists.

CLT

Theorem: Lindeberg-Levy CLT:

$\{X_j\}$ iid r.v. $E[X_j] = \mu$, $\text{Var}(X_j) = \sigma^2$, then: $\sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \xrightarrow{d} N(0, 1)$.

Properties of Normal Distribution:

$$X_i \sim N(\mu, \sigma^2) \quad \sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2) \quad \frac{1}{n} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} \sim N(0, 1) \quad \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} \right) \sim N(0, 1)$$

Theorem: Cramér-Wold device

$\{X_n\}$ a sequence of random k-vectors. X random k-vector. Then the followings are equivalent

- ① $X_n \xrightarrow{d} X$ ② $\lambda' X_n \xrightarrow{d} \lambda' X$ for all $\lambda \in \mathbb{R}^k$.

Theorem: Multivariate CLT:

$\{X_i\}$: iid random k-vectors. $E[X]$: common mean of X_i . $V = \text{Var}(X_i)$ common covariance matrix (finite non-singular)
 $\sqrt{n} (\bar{X}_n - E[X]) \xrightarrow{d} N(0, V)$.

The Delta Method.

Theorem: Univariate delta method:

$\{z_n\}$ r.v. $g: \mathbb{R} \rightarrow \mathbb{R}$. Let $\mu \in \mathbb{R}$.

Suppose that $g'(\mu)$ exists, $\neq 0$, continuous at μ .

Suppose: $\sqrt{n}(z_n - \mu) \xrightarrow{d} N(0, \sigma^2)$ as $n \rightarrow \infty$

then $\sqrt{n}(g(z_n) - g(\mu)) \xrightarrow{d} N(0, g'(\mu)^2 \sigma^2)$

Pf: Taylor's expansion near μ :

$$g(z_n) = g(\mu) + g'(\tilde{z}_n)(z_n - \mu), \quad \tilde{z}_n \text{ between } \mu \text{ and } z_n$$

$$g(z_n) - g(\mu) = g'(\tilde{z}_n)(z_n - \mu) \quad \tilde{z}_n \xrightarrow{P} \mu.$$

$$\sqrt{n}(g(z_n) - g(\mu)) = \sqrt{n} g'(\tilde{z}_n)(z_n - \mu)$$

$$\therefore \sqrt{n}(z_n - \mu) \xrightarrow{d} N(0, \sigma^2) \therefore z_n \rightarrow \mu \therefore \tilde{z}_n \rightarrow \mu.$$

$$\therefore \sqrt{n}(g(z_n) - g(\mu)) \xrightarrow{d} N(0, g'(\mu)^2 \sigma^2).$$

Theorem: Multivariate Delta Method.

$g: \mathbb{R}^k \rightarrow \mathbb{R}^l$ have non-zero derivative $\nabla g(\alpha)$ at $\alpha \in \mathbb{R}^k$.

$b > 0$ scalar. $\{z_n\}$ r. k-vectors. \bar{z} r. k-vector

Suppose: $n^b (z_n - \alpha) \xrightarrow{d} \bar{z}$

then: $n^b (g(z_n) - g(\alpha)) \xrightarrow{d} [\nabla g(\alpha)]^2 \bar{z}$

e.g. $\{x_i\}$ iid $f(x) = \alpha e^{-\alpha x}$. $E[x] = \frac{1}{\alpha}$ $\alpha = \frac{1}{E[x]}$

$$\hat{\alpha} = \frac{1}{\bar{x}_n} \quad \text{By CLT: } \hat{\alpha} \xrightarrow{P} \alpha. \quad \text{Var}(\bar{x}) = \frac{1}{n}$$

$$\therefore g(\bar{x}_n) = g\left(\frac{1}{\bar{x}_n}\right) = \alpha \quad g'(\bar{x}_n) = g'\left(\frac{1}{\bar{x}_n}\right) = -\frac{1}{\bar{x}_n^2}$$

$$\therefore g(E[x]) = g\left(\frac{1}{\alpha}\right) = \alpha \quad g'(E[x]) = g'\left(\frac{1}{\alpha}\right) = -\frac{1}{(\alpha)^2}$$

$$\therefore \sqrt{n}(\alpha - \hat{\alpha}) \xrightarrow{d} N(0, [-\frac{1}{(\alpha)^2}]^2 \cdot \frac{1}{n})$$

$$\therefore \hat{\alpha} \xrightarrow{d} N\left(\alpha, \frac{\alpha^2}{n}\right)$$

Theory of Statistics.

A: action space $\mathcal{P} \in \mathcal{P}$: state of the world. $U(P, a)$: utility function $X = (X_1 \dots X_n)$: observed realization

$d(X) \in A$: our decision

Def. $d: \mathbb{R}^n \rightarrow A$: a statistical decision rule.

Want to solve: $\max_{d(X)} U(P, d(X))$

Subjective Bayesian Approach: ex post

① have prior belief $\pi(\cdot)$ about P .

② update it to posterior belief $\pi(\cdot | x)$ after seeing x

③ $\max_{d(x)} \int U(P, d(x)) d\pi(P|x)$.

Problems: ① $\pi(\cdot)$ come from?

② P may be non-parametric.

Frequentist's Approach: ex ante

Def: the distribution of d is the sampling distribution,

where $d: \mathbb{R}^n \rightarrow A$. (dist. of $(X_1 \dots X_n)$) induces it via the r.v. $d(X)$
depends on P and d .

Def. welfare function: (ex ante) $W(P, d) = \int U(P, d(x)) dP(x)$
unknown since P unknown.

$$= \mathbb{E}_P[U(P, d(x))]$$

eliminate X .

2 ways to get rid of P :

①

(i) suppose a prior dist. $\pi(\cdot)$ of P : $B(d) = \int W(P, d) d\pi(P)$

(ii) $\max_{d(\cdot)} B(d) \rightarrow$ Bayesian optimal decision rule. d^*

②. Worst case scenario:

(i) for each choice of d : $\text{Worstcase}(d) = \min_{P \in \mathcal{P}} W(P, d)$

(ii). $\max_{d \in \mathcal{D}} \text{Worstcase}(d) \rightarrow$ Maximin optimal decision rule d^* .

Point Estimation.

P : probability measure of interest. $\theta = \theta(P)$: parameter of interest. \mathbb{H} : parameter space (set of all logically possible values).

Def. $X_1 \dots X_n$ iid. A point estimator is a measurable function $\hat{\theta}: \mathbb{R}^n \rightarrow \mathbb{H}$.

$(\hat{\theta}) = \hat{\theta}(X_1 \dots X_n)$: estimator $\hat{\theta}(x_1 \dots x_n)$: estimate (with realized $x_1 \dots x_n$).

Good $\hat{\theta}(\cdot)$: ① close ② precise.

Finite Sample Version:

Closedness: $\text{Bias}(\hat{\theta}, P) = \mathbb{E}[\hat{\theta}] - \theta = \mathbb{E}_P(\hat{\theta}(X_1 \dots X_n)) - \theta$
 $= 0$, then unbiased.

Precision: $\text{Var}(\hat{\theta}) = \text{Var}_P(\theta(X_1 \dots X_n))$ as small as possible.

Balancing: $\text{MSE}(\hat{\theta}, P) = \mathbb{E}_P[(\hat{\theta}(x) - \theta)^2] = \text{Bias}^2(\hat{\theta}, P) + \text{Var}(\hat{\theta}, P)$.
 if $\text{MSE}(\hat{\theta}_1, P) \leq \text{MSE}(\hat{\theta}_2, P)$ $\forall n \in \mathbb{N}$, then $\hat{\theta}_1$ is relatively more efficient than $\hat{\theta}_2$.

Minimax optimal estimator:

$$\min_{\hat{\theta}(\cdot)} \max_{P \in \mathcal{P}} \text{MSE}(\hat{\theta}, P).$$

Asymptotic Version:

Closedness:

1. (Weakly) consistency: $\hat{\theta} \xrightarrow{P} \theta$ for all P , $\theta(P)$ as $n \rightarrow \infty$.

2. Unbiased in the limit: $\mathbb{E}[\hat{\theta}] \rightarrow \theta$ as $n \rightarrow \infty$.

3. Asymptotically unbiased: $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} Z$ if $\mathbb{E}[Z] = 0$, $\sqrt{n} \rightarrow \infty$. Z non-degenerate r.v.
 $3 \rightarrow 1$

Precision:

$\text{Var}(\hat{\theta}) \rightarrow 0$ in most cases. So, we look at:

1. $\lim_{n \rightarrow \infty} \text{Var}(\sqrt{n}(\hat{\theta} - \theta))$, where $\sqrt{n}(\hat{\theta} - \theta) \rightarrow Z$. $\left\{ \begin{array}{l} \sqrt{n}(\hat{\theta}_1 - \theta) \rightarrow Z_1 \\ \sqrt{n}(\hat{\theta}_2 - \theta) \rightarrow Z_2 \end{array} \right. \text{e.g. } \text{var}(Z_1) > \text{var}(Z_2) \text{ then } \hat{\theta}_2 \text{ more precise.}$
2. $\sqrt{n}(\hat{\theta} - \theta) \rightarrow Z$, look at the dist. of Z , e.g. $\text{var}(Z)$.

Hypothesis Tests

$X_1 \dots X_n$: iid P , $\theta = \theta(P)$. $\mathbb{H} = \mathbb{H}_0 \cup (\mathbb{H} \setminus \mathbb{H}_0)$ $\theta \in \mathbb{H}_0$ or $\theta \in \mathbb{H} \setminus \mathbb{H}_0$.

0: accepting \mathbb{H}_0 .

Definition: A hypothesis test is a measurable function $\phi_n = \phi_n(X_1 \dots X_n): \mathbb{R}^n \rightarrow \{0, 1\}$. 1: rejecting \mathbb{H}_0 .

Def. Power Function: Let ϕ_n be a test. Define the function $\text{power}: \mathcal{P} \rightarrow [0, 1]$ by: $\text{Power}(P, \phi_n) = P(\phi_n(X) = 1)$.

Def. Risk function of the test ϕ_n is:

$$\begin{aligned} R(P, \phi_n) &= L(P, 1)P(\phi_n(X) = 1) + L(P, 0)P(\phi_n(X) = 0) \\ &= L(P, 0) + P(\phi_n(X) = 1)(L(P, 1) - L(P, 0)). \end{aligned}$$

Where $L(P, a)$ is the loss function: $\phi_n(x) \rightarrow P(\phi_n(x) = 1) \rightarrow R(P, \phi_n)$

Def. 0-1 loss function (most standard one): $L(P, a) = \begin{cases} 1 & \text{if } \theta(P) \in \mathbb{H}_0 \text{ and } a = 1 \\ 0 & \text{otherwise.} \end{cases}$

False positive

False negative

Under 0-1 loss function, $R(\bar{P}, \phi_n) = \begin{cases} \bar{P}(\phi_n(x)=1) & \text{if } \Theta(\bar{P}) \in \mathbb{H}_0 \\ -\bar{P}(\phi_n(x)=1) & \text{if } \Theta(\bar{P}) \notin \mathbb{H}_0 \end{cases}$

Best power function: $\text{Power}(\bar{P}, \phi_n^*) = \begin{cases} 1 & \text{if } \Theta(\bar{P}) \notin \mathbb{H}_0 \\ 0 & \text{if } \Theta(\bar{P}) \in \mathbb{H}_0 \end{cases}$ not feasible! Instead, find "good one" in some sense.

The Neyman-Pearson Paradigm.

First, restrict attention to tests whose prob. of false positive is not too big.

Def. Size of a test ϕ_n : $\text{size}(\phi_n) = \sup_{\bar{P} \in \mathcal{P}: \Theta(\bar{P}) \in \mathbb{H}_0} \bar{P}(\phi_n(x)=1)$, (largest prob. of false positive under all $\bar{P} \in \mathcal{P}$).

Def. ϕ_n is a level α test if $\text{size}(\phi_n) \leq \alpha$, $\alpha \in (0, 1)$

Next, choose the one that has the largest power.

Def: ϕ_n^* is the uniformly most powerful (UMP) within \mathbb{B} if, for any $\phi_n \in \mathbb{B}$,

$$\text{Power}(\bar{P}, \phi_n^*) \geq \text{Power}(\bar{P}, \phi_n)$$

for all \bar{P} such that $\Theta(\bar{P}) \notin \mathbb{H}_0$.

Problem & solution: see weekly summary.

Cutoff test. (Most hypothesis tests have this form):

$$\phi_n(x) = \mathbb{I}[T(x) > c]$$

$T(x)$: test statistic. c : critical value.

Asymptotic Criteria:

Def. the uniform asymptotic size of ϕ_n is: $\text{Asympsize}(\phi_n) = \limsup_{n \rightarrow \infty} \sup_{\bar{P} \in \mathcal{P}: \Theta(\bar{P}) \in \mathbb{H}_0} \bar{P}(\phi_n(x)=1)$.

Def. ϕ_n is pointwise consistent if $\lim_{n \rightarrow \infty} \bar{P}(\phi_n(x)=1) = 1$ for all $\bar{P} \in \mathcal{P}$ such that $\Theta(\bar{P}) \notin \mathbb{H}_0$.

P-values.

Def. let \mathbb{B} be the set of all non-randomized tests of the null hypothesis \mathbb{H}_0 . Let $x \in \text{supp}(X)$.

$$p(\mathbb{H}_0, x) = \inf_{\phi_n \in \mathbb{B}: \phi_n(x)=1} \sup_{\bar{P} \in \mathcal{P}: \Theta(\bar{P}) \in \mathbb{H}_0} \bar{P}(\phi_n(x)=1)$$

$$= \inf_{\phi_n \in \mathbb{B}: \phi_n(x)=1} \text{size}(\phi_n)$$

The definition based on a specific cutoff test:

Let the test be $\phi_n(x) = \mathbb{I}[T(x) > c]$

Let $c = c(1-\alpha)$. $c'(\cdot) > 0$. As $\alpha \uparrow$, $c \downarrow$. As $\alpha \downarrow$, $c \uparrow$. (size \downarrow , critical value \uparrow , harder to reject).

$$\phi_n(x) = \mathbb{I}[T(x) > c(1-\alpha)]$$

Let $x \in \text{supp}(X)$ be a specific realization of the data, so $T(x)$ is fixed.

Suppose that we reject. $T(x) > c(1-\alpha)$.

This is not strict enough since $T(x) - c(1-\alpha) > 0$, so we let $\alpha \downarrow$ to p so that $c(1-\alpha) \uparrow$ to $c(1-p)$

$$T(x) = c(1-p)$$

$p = 1 - \mathbb{P}(T(x))$: this is the p-value within the restrict class of \mathbb{B} .

①. $p < \alpha$.

②. as $\alpha \downarrow$ to p , still can reject \mathbb{H}_0 because $T(x) = c(1-p) > c(1-\alpha)$

③. as $\alpha \downarrow$ to be less than p , cannot reject \mathbb{H}_0 because $T(x) = c(1-p) < c(1-\alpha)$ now.

Add an assumption : $C^{-1}(t) = \bar{P}_0(T(X) \leq t)$ *

\bar{P}_0 : the dist. of X that $\Theta(\bar{P}_0) = \Theta$. $C^{-1}(t)$: cdf of $T(x)$ when true dist. is \bar{P}_0 . $(\cdot)^{-1}$: inverse of cdf: quantile function.

Then: $p(x) = 1 - \bar{P}_0(T(X) \leq T(x))$ Most common def: p-value for a cutoff test with $T(\cdot)$ and *
 $= \bar{P}_0(T(X) > T(x))$. is the prob. of observing a test statistic's value more extreme than the one we actually saw in the data, assuming that the true dist is \bar{P}_0 .

p-value depend on x , so it is a r.v.

Lemma. define $p(x) = \bar{P}_0(T(X) \geq T(x))$.

Suppose $X \sim \bar{P}_0$. Then $p \sim \text{Unif}[0, 1]$. (ex-ante).

$$\text{Pf: } \bar{P}_0(p(x) \leq p) = \bar{P}_0(1 - F_{T(X)}(T(x)) \leq p)$$

$$= \bar{P}_0(1-p \leq F_{T(X)}(T(x)))$$

$$= \bar{P}_0(F_{T(X)}^{-1}(1-p) \leq T(x))$$

$$= 1 - F_{T(X)}(F_{T(X)}^{-1}(1-p))$$

$$= 1 - (1-p)$$

$$= p.$$

Theorem. define $\phi_n(x) = \bar{P}_0(T(X) > T(x))$. Let $\alpha \in [0, 1]$.

define: $\phi_n(x) = \mathbb{I}[\phi_n(x) \leq \alpha]$

Then, $\phi_n(x)$ is a level- α test of the simple null \bar{P}_0

against the alternative that $P \neq \bar{P}_0$.

$$\text{Pf: } \bar{P}_0(\phi_n(x) = 1) = \bar{P}_0(p(x) \leq \alpha)$$

$= \alpha$ according to the Lemma.

↑ using $p(x)$ to construct level- α test.

Confidence Set.

it depends on sample size!

Def. A level α confidence set for Θ is a function $C(X) = C(X_1, \dots, X_n)$ from the sample to the subset of Θ such that 1. $\{x \in \text{supp}(X) : C(x) \ni \Theta(P)\}$ is P -measurable and 2. $\bar{P}(C(x) \ni \Theta(P)) \geq 1-\alpha$ for all $P \in P$.

Interpretation: $\Theta(P)$ is fixed. throw sets at it and hit it at least $100(1-\alpha)\%$ of the times.

Asymptotic version: $\liminf_{n \rightarrow \infty} \inf_{P \in P} \bar{P}(C(x) \ni \Theta(P)) \geq 1-\alpha$.

Rejection region $R = \{x \in \text{supp}(X) : \phi_{\theta_0}(x) = 1\}$ flip the
 Tolerance confidence set $C(x) = \{\theta \in \Theta : \phi_{\theta}(x) = 0\}$ inequality!

- ① For each $\theta \in \Theta$, let $\phi_{\theta}(x)$ be a level α test of $H_0: \Theta(P) = \theta$. Let $C(x) = \{\theta \in \Theta : \phi_{\theta}(x) = 0\}$, then $C(x)$ is a level α confidence set.
- ② Let $C(x)$ be a level α confidence set. For each $\theta \in \Theta$, define $\phi_{\theta}(x) = \mathbb{I}[C(x) \ni \theta]$, then for each $\theta \in \Theta$, $\phi_{\theta}(x)$ is a level α test of $H_0: \Theta(P) = \theta$.

Regression Mechanism.

describing the conditional mean of an observed r.v. $Y | X=x$. Study $E[Y | X=x]$

Basic Variables X

$$\text{Derived Variables } p(x) = \begin{bmatrix} p_1(x) \\ \vdots \\ p_j(x) \end{bmatrix}$$

proof:

$$E[Y | X=x] = p(x)' \beta$$

$$p(x) E[Y | X=x] = p(x) p(x)' \beta$$

$$E[p(x) E[Y | X=x]] = E[p(x) p(x)'] \beta$$

$$E[p(x) Y] = E[p(x) p(x)'] \beta$$

$$\beta = E[p(x) p(x)']^{-1} E[p(x) Y]$$

Suppose $m(x) = E[Y | X=x] = p(x)' \beta$, then.

Theorem: Y : r.v. X : r.vector. Assume:

① linearity: $m(x) = E[Y | X=x] = p(x)' \beta$

② finite moments: $E[p(x)Y]$ and $E[p(x)p(x)']$ finite

③ sufficient variation: $E[p(x)p(x)']$ invertible.

$$\text{Then: } \beta = E[p(x)p(x)']^{-1} E[p(x)Y]$$

$$\text{proof: } p(x)p(x)' = \begin{bmatrix} 1 & | & x \\ x & | & \end{bmatrix} \begin{bmatrix} 1 & | & x \\ x & | & \end{bmatrix}' = \begin{bmatrix} 1 & | & x \\ x & | & x^2 \end{bmatrix} E[p(x)p(x)'] = \begin{bmatrix} 1 & | & E[x] \\ E[x] & | & E[x^2] \end{bmatrix}$$

$$E[p(x)p(x)']^{-1} = \frac{1}{E[x^2] - E[x]^2} \begin{bmatrix} E[x] - E[x] \\ -E[x] & 1 \end{bmatrix}$$

$$E[p(x)p(x)']^{-1} E[p(x)Y] = \frac{1}{E[x^2] - E[x]^2} \begin{bmatrix} E[x] - E[x] \\ -E[x] & 1 \end{bmatrix} \begin{bmatrix} E[Y] \\ E[XY] \end{bmatrix}$$

$$= \frac{1}{E[x^2] - E[x]^2} \begin{bmatrix} E[x]E[Y] - E[X]E[XY] \\ E[XY] - E[X]E[Y] \end{bmatrix}$$

Proposition: Let Y and X be random variables. $p = (1, x)'$.

$$\beta^{obs} = \begin{pmatrix} \beta_0^{obs} \\ \beta_1^{obs} \end{pmatrix} = E[p(x)p(x)']^{-1} E[p(x)Y]$$

$$\text{Then: } \beta_1^{obs} = \frac{Cov(Y, X)}{\text{var}(X)} \cdot \beta_0^{obs} = E[Y] - \beta_0^{obs} E[X].$$

Sufficient variation: $\text{var}(X) > 0$.

The residual zero correlation property.

Proposition: Define the residual. $E \equiv Y - p(x)' \beta^{obs}$, then.

$$\textcircled{1} \quad \mathbb{E}[p(x)E] = 0.$$

If $p(x)$ contains a constant, then $\mathbb{E}(E) = 0$, $\text{Cov}(p_j(x), E) = 0 \cdot j=1, 2, \dots, J$

$$\textcircled{2} \quad \text{Cov}(Y, E) = \text{Var}(E).$$

proof: $\textcircled{1} \quad \mathbb{E}[p(x)E] = \mathbb{E}[p(x)(Y - p(x)\beta^{ou})]$
 $= \mathbb{E}[p(x)Y] - \mathbb{E}[p(x)p(x)\beta^{ou}]$
 $= \mathbb{E}[p(x)Y] - \mathbb{E}[p(x)p(x)]\mathbb{E}[p(x)p(x)]^T \mathbb{E}[p(x)Y]$
 $= 0$

$$\textcircled{2} \quad \text{Cov}(Y, E) = \text{Cov}(p(x)\beta^{ou} + E, E)$$

 $= \sum_{j=1}^J \text{Cov}(p_j(x), E)\beta_j^{ou} + \text{Cov}(E, E)$
 $= \text{Var}(E)$

Residual Regression: the FWL theorem.

$$\text{Let } X_{-k} = (1, X_1, \dots, X_{k-1}, X_{k+1}, \dots, X_k)^T$$

$$\text{let } \tilde{\beta}_k = \mathbb{E}[X_k X_k^T]^{-1} \mathbb{E}[X_k Y], \quad X_k^{\perp X_k} = X_k - X_k \tilde{\beta}_k$$

Theorem: FWL theorem: Let Y be r. v., X r. K -vector.

$$\text{Let } p(x) = (1, x)^T. \text{ Define } \beta^{ou} = \mathbb{E}[p(x)p(x)]^{-1} \mathbb{E}[p(x)Y].$$

$$\text{Then: } \beta_k^{ou} = \frac{\text{Cov}(Y, X_k^{\perp X_k})}{\text{Var}(X_k^{\perp X_k})}$$

proof: define $E = Y - p(x)\beta^{ou}$

$$\begin{aligned} \text{Cov}(Y, X_k^{\perp X_k}) &= \text{Cov}(\beta_0^{ou} + \beta_1 X_1 + \dots + \beta_K X_K + E, X_k^{\perp X_k}) \\ &= \beta_0^{ou} \text{Cov}(X_1, X_k^{\perp X_k}) + \beta_1 \text{Cov}(X_2, X_k^{\perp X_k}) + \dots + \text{Cov}(E, X_k^{\perp X_k}) \\ \text{Cov}(E, X_k^{\perp X_k}) &= \text{Cov}(E, X_k - \sum_{j \neq k} X_j \tilde{\beta}_{kj}) \\ &= \text{Cov}(E, X_k) - \sum_{j \neq k} \tilde{\beta}_{kj} \text{Cov}(E, X_j) \\ &= 0 \end{aligned}$$

For $j \neq k$: $\text{Cov}(X_j, X_k^{\perp X_k}) = 0$ follows $\textcircled{1}$ of proposition

For $j = k$: $\text{Cov}(X_k, X_k^{\perp X_k}) = \text{Var}(X_k^{\perp X_k})$ $\textcircled{2}$ of proposition.

$$\therefore \beta_k^{ou} = \frac{\text{Cov}(Y, X_k^{\perp X_k})}{\text{Var}(X_k^{\perp X_k})}$$

The long and short regressions.

Analyze why and when coefficients change when adding variables.

Short regression: Y on $(1, x)$; get $(\alpha^{short}, \beta^{short})$

Long regression: Y on $(1, x, w)$; get $(\alpha^{long}, \beta^{long}, \gamma^{long})$

$$\text{Then: } \beta^{short} = \frac{\text{Cov}(Y, x)}{\text{Var}(x)} \quad \beta^{long} = \frac{\text{Cov}(Y, x^w)}{\text{Var}(x^w)} \quad \beta^{short} = \beta^{long} + \gamma^{long} \frac{\text{Cov}(w, x)}{\text{Var}(x)}$$

proof. define $E = Y - (\alpha^{long} + \beta^{long}x + \gamma^{long}w)$

$$\begin{aligned} \text{Then: } \text{Cov}(Y, x) &= \text{Cov}(\alpha^{long} + \beta^{long}x + \gamma^{long}w + E, x) \\ &= \beta^{long} \text{Var}(x) + \gamma^{long} \text{Cov}(w, x) + \text{Cov}(E, x) \\ &= \beta^{long} \text{Var}(x) + \gamma^{long} \text{Cov}(w, x). \quad \text{follows } \textcircled{1} \text{ of proposition.} \end{aligned}$$

$$\therefore \frac{\text{Cov}(Y, x)}{\text{Var}(x)} = \beta^{short} = \beta^{long} + \gamma^{long} \frac{\text{Cov}(w, x)}{\text{Var}(x)}.$$

Characterize OLS coefficients.

$\textcircled{1}$ Apply FWL theorem

$\textcircled{2}$ show that the conditional mean function is linear.

Saturated regression:

$$m(x) = \mathbb{E}[Y | X=x] = p(x)\beta$$

If we pick the right $p(x)$, the linearity is guaranteed to hold when X is discrete.

$$\text{Let } X = \{0, 1\}, \quad p(x) = (1, \mathbf{1}(x=1))^T$$

$$\text{Define } \beta_0 = \mathbb{E}[Y | X=0], \quad \beta_1 = \mathbb{E}[Y | X=1] - \mathbb{E}[Y | X=0]$$

$$\text{Then: } m(x) = \mathbb{E}[Y | X=x] = p(x)\beta = \mathbb{E}[Y | X=0] + \mathbf{1}(x=1)(\mathbb{E}[Y | X=1] - \mathbb{E}[Y | X=0]) \\ = \beta_0 + \beta_1 x$$

Regression Reinterpretation: interpreting OLS estimand when the true conditional mean function is not linear.

As a best linear approximation to the true non-linear conditional mean function

As a best linear predictor of the outcome variable

As a weighted average derivative of the true non-linear conditional mean function.

OLS estimator

Ways to come up with estimator:

① Analogy principle: want to know $\theta = \mathbb{E}[F_X]$, use the sample colf to estimate $\hat{\theta} = \mathbb{E}[\hat{F}_X]$

$$F_X(x) = \mathbb{E}[\mathbb{I}[X \leq x]], \hat{F}_X(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[X_i \leq x]. \quad \theta = \mathbb{E}[X], \quad \hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$$

② MLE ③ GMM

Obtaining OLS as a sample analogue:

$$\beta = \mathbb{E}[XX'] \mathbb{E}[XY]$$

$$\hat{\beta}_{OLS} = \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i Y_i \right)$$

Obtaining OLS via least squares:

$$\beta = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} \mathbb{E}[(Y - X'b)^2]$$

$$\hat{\beta}_{OLS} = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - X_i b)^2.$$

Observe: $\{(y_i, x_i)\}_{i=1}^n$
 where $x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ik} \end{bmatrix}$ $X = \begin{bmatrix} x_1' \\ \vdots \\ x_n' \end{bmatrix} = \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nk} \end{bmatrix}$ $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

$$\text{For each } i: e_i = y_i - x_i' \beta \quad y_i = x_i' \beta + e_i$$

$$\text{Stack up from } i=1 \text{ to } n: \quad y = X \beta + e \quad e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

For $\beta = \mathbb{E}[XX']^{-1} \mathbb{E}[XY]$:

① Asymptotic variance matrix without hetero:
 $V_1 = \sigma^2 \mathbb{E}[XX']^{-1}$

② with hetero: $V_2 = \mathbb{E}[XX']^{-1} \mathbb{E}[E^2 XX'] \mathbb{E}[XX']^{-1}$

$$\text{for an } n \times 1 \text{ matrix } a: a'a = \sum_{i=1}^n a_i^2$$

$$\text{for a matrix } X: X'X = [x_1 \dots x_n] \begin{bmatrix} x_1' \\ \vdots \\ x_n' \end{bmatrix}$$

$$= \sum_{i=1}^n x_i x_i'$$

$$\rightarrow \text{Therefore, equivalent as minimizing } \sum_{i=1}^n (y_i - x_i' \beta)^2 = (y - X\beta)'(y - X\beta) = y'y - b'X'y - y'Xb + b'X'Xb$$

$$\therefore \frac{d \sum_{i=1}^n (y_i - x_i' \beta)^2}{d \beta} = -X'y - X'b + 2X'Xb = 0$$

$$\therefore \hat{\beta}_{OLS} = (X'X)^{-1} X'y \quad (\hat{\beta})$$

$$\frac{d(a'b)}{db} = \frac{d(b'a)}{db} = a$$

$$\frac{d(b'A'b)}{db} = (A+A')b = 2Ab \quad \text{if } A \text{ symmetric.}$$

$$\hat{\beta}_{OLS} = (X'X)^{-1} (X'Y)$$

$$= \left(\sum_{i=1}^n x_i x_i' \right)^{-1} \left(\sum_{i=1}^n x_i y_i \right) = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right) \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) \quad \text{Same as analogue!}$$

Finite sample residuals.

$$-X'y + X'Xb = 0 \Rightarrow X'(y - Xb) = 0$$

define finite sample residual for the i th observation by: $\hat{e}_i = y_i - x_i' \hat{\beta}$ (an estimation of population level residual $e_i = y_i - x_i' \beta$)

Stacking up: $\hat{e} = y - X\hat{\beta}$

then: $X' \hat{e} = \begin{bmatrix} x_1' \\ \vdots \\ x_n' \end{bmatrix}' \hat{e} = (x_1 \dots x_n)' \hat{e} = 0$, a system of $(x_{1k} \dots x_{nk})' \hat{e} = \langle (x_{1k} \dots x_{nk}), \hat{e} \rangle = 0$

$\therefore \hat{e} \perp X$, $\hat{e} \perp$ every column of X ($(x_{1k} \dots x_{nk})' \perp \hat{e}$)

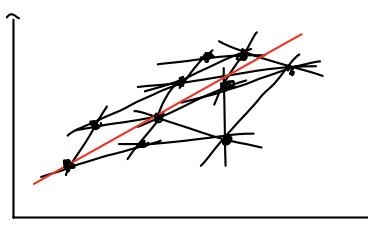
(population: $E \perp X$, $\mathbb{E}[X_k E] = 0 \Leftrightarrow \frac{1}{n} \sum_{i=1}^n x_{ik} \hat{e}_i = 0$. When X_i contains a constant, $\frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0 \Leftrightarrow \mathbb{E}[\hat{e}] = 0$)

Geometric interpretation of OLS

Obtaining OLS as an MLE.

Assume $Y|X=x \sim N(X'\beta, \sigma^2)$, then $\hat{\beta}_{OLS} = \hat{\beta}_{MLE}$

OLS as a weighted average of pairwise slope.



OLS (finite sample)

Unbiasedness:

Proposition: Unbiased under linearity.

If $\mathbb{E}[Y|X=x] = x'\beta$, $\{(y_i, x_i)\}_{i=1}^n$ iid. from $(Y|X)$
then $\mathbb{E}[\hat{\beta}_{OLS}] = \beta$

$$\begin{aligned}\text{proof: } \mathbb{E}[\hat{\beta}_{OLS}|X] &= \mathbb{E}[(X'X)^{-1}X'y|X] \rightarrow \mathbb{E}[y_i|x_1 \dots x_n] \\ &= (X'X)^{-1}X'\mathbb{E}[y|X] \\ &= (X'X)^{-1}X'\times\beta \\ &= \beta \\ \mathbb{E}[\hat{\beta}_{OLS}] &= \mathbb{E}[\mathbb{E}[\hat{\beta}_{OLS}|X]] = \beta\end{aligned}$$

Variance of OLS:

Def. Conditional variance function: $\sigma^2(x) = \text{Var}(Y|X=x)$

Proposition: Assume: ① linear $\mathbb{E}[Y|X=x] = x'\beta$ ② iid $\{(y_i, x_i)\}_{i=1}^n$

$$\text{then: } \text{Var}(\hat{\beta}_{OLS}|X) = (X'X)^{-1}X'\Omega X(X'X)^{-1}$$

$$\Omega = \begin{bmatrix} \sigma^2(x_1) & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & \sigma^2(x_n) \end{bmatrix} = \text{Var}(y|X)$$

If homoskedasticity, $\text{Var}(\hat{\beta}_{OLS}|X) = \sigma^2(X'X)^{-1}$ ($V = \sigma^2 \mathbb{E}[xx']^{-1}$)

For either case, $\text{Var}(\hat{\beta}_{OLS}) = \mathbb{E}[\text{Var}(\hat{\beta}_{OLS}|X)]$

The law of iterated variance: $\text{Var}(Y) = \text{Var}(\mathbb{E}[Y|X]) + \mathbb{E}[\text{Var}(Y|X)]$

$$\begin{aligned}\therefore \text{Var}(\hat{\beta}_{OLS}) &= \mathbb{E}[\text{Var}(\hat{\beta}_{OLS}|X)] + \text{Var}(\mathbb{E}[\hat{\beta}_{OLS}|X]) \\ &= \mathbb{E}[\text{Var}(\hat{\beta}_{OLS}|X)] \quad \downarrow \text{constant}\end{aligned}$$

Matrix Analogue:

$$\mathbb{E}[A \cdot y] = A \cdot \mathbb{E}[y]$$

$$\text{Var}(Ay) = A \cdot \text{Var}(y) \cdot A'$$

$$\text{Var}(y) = \begin{bmatrix} \text{Var}(y_1) & \text{Cov}(y_1, y_2) & \cdots & \text{Cov}(y_1, y_n) \\ \text{Cov}(y_2, y_1) & \text{Var}(y_2) & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(y_n, y_1) & \text{Cov}(y_n, y_2) & \cdots & \text{Var}(y_n) \end{bmatrix}$$

$$\begin{aligned}\text{proof: } \text{Var}(\hat{\beta}|X) &= \text{Var}((X'X)^{-1}X'y|X) \\ &= (X'X)^{-1}X'\text{Var}(y|X)X(X'X)^{-1}\end{aligned}$$

for $\text{Var}(y|X)$: its diagonal entries are:

$$\begin{aligned}\text{Var}(y_i|X) &= \text{Var}(y_i|x_1 \dots x_n) \\ &= \text{Var}(y_i|x_i) \\ &= \sigma^2(x_i)\end{aligned}$$

for off-diagonal entries:

$$\text{Cov}(y_i, y_j|X) = \mathbb{E}[y_i y_j|X] - \mathbb{E}[y_i|X] \mathbb{E}[y_j|X] \quad (\star)$$

$$\mathbb{E}[y_i y_j|X] = \mathbb{E}[y_i y_j|x_1 \dots x_n]$$

$$= \mathbb{E}[y_i y_j|x_i, x_j]$$

$$= \mathbb{E}[\mathbb{E}[y_i y_j|x_i, x_j, y_{\bar{i}, \bar{j}}]|x_i, x_j]$$

$$= \mathbb{E}[y_i \mathbb{E}[y_j|x_i, x_j, y_{\bar{i}, \bar{j}}]|x_i, x_j]$$

$$= \mathbb{E}[y_i|x_i] \mathbb{E}[y_j|x_i, x_j]$$

$$= \mathbb{E}[y_i|x_i] \mathbb{E}[y_j|x_j]$$

$$\therefore (\star) = \mathbb{E}[y_i|x_i] \mathbb{E}[y_j|x_j] - \mathbb{E}[y_i|x_i] \mathbb{E}[y_j|x_j]$$

$$= 0$$

The Variance in the scalar regressor case:

OLS of Y on X : design matrix is: $X = (1, x)$, 1 is a column vector of ones.

$$X'X = \begin{bmatrix} 1' \\ x' \end{bmatrix} \begin{bmatrix} 1 & x \end{bmatrix} = \begin{bmatrix} 1'1 & 1'x \\ x'1 & x'x \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$\therefore (X'X)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

Suppose homoskedasticity. then

$$\text{Var}(\hat{\beta}_1 | X) = \frac{\sigma^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

① $\sigma^2 \rightarrow 0 : \text{Var}(\hat{\beta}_{OLS}|X) \rightarrow 0$.

② $n \rightarrow \infty : \text{Var}(\hat{\beta}_{OLS}|X) \rightarrow 0$.

③ $\text{Var}(x) \rightarrow \infty : \text{Var}(\hat{\beta}_{OLS}|X) \rightarrow 0$.

$$\begin{aligned}\sum_i (x_i - \bar{x})^2 &= \sum_i (x_i^2 - 2x_i \bar{x} + \bar{x}^2) \\ &= \sum_i x_i^2 - \frac{2}{n} \sum_i x_i \cdot \sum_i x_i + \frac{1}{n} (\sum_i x_i)^2 \\ &= \sum_i x_i^2 - \frac{1}{n} (\sum_i x_i)^2\end{aligned}$$

$$\begin{aligned}\text{Var}(\hat{\beta}_1 | X) &= \frac{n \sigma^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} = \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \\ &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{n \text{Var}(x)}\end{aligned}$$

OLS is the best linear unbiased estimator of β (BLUE).

Estimating σ^2 to do inference on $\hat{\beta}$

$$\sigma^2 = \text{Var}(Y) = \text{Var}(E) = \mathbb{E}[E^2] - (\mathbb{E}[E])^2 = \mathbb{E}[E^2] \quad (\text{homoskedasticity}) \quad \text{Var}(\hat{\beta}|X) = \sigma^2 (X'X)^{-1}$$

population: $E = Y - X'\beta$

infeasible: $e_i = y_i - x_i'\beta = y_i - x_i'\hat{\beta}$ $\hat{e} = \begin{bmatrix} y_1 - x_1'\hat{\beta} \\ \vdots \\ y_n - x_n'\hat{\beta} \end{bmatrix}$

feasible: $\hat{e}_i = y_i - x_i'\hat{\beta}_{OLS}$

$$\therefore \text{from sample analogue: } \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2 = \frac{1}{n} \hat{e}' \hat{e} \quad (\text{biased}).$$

proposition: Assume: ① linearity $E[Y|X=x] = X'\beta$ ② iid $\{y_i, x_i\}_{i=1}^n$ ③ $\sigma^2(x) = \text{Var}(Y|X=x) = \sigma^2 > 0$

$$\text{Then: } E[\hat{\sigma}^2|X] = E\left[\frac{1}{n}\hat{e}'\hat{e}|X\right] = \frac{n-k}{n}\sigma^2 \quad E[\hat{\sigma}^2] = \frac{n-k}{n}\sigma^2 \text{ Biased.}$$

proof: $M = I_n - P$ where $P = X(X'X)^{-1}X'$

$$\begin{aligned} \hat{e} &= y - X\hat{\beta} & \hat{e}'\hat{e} &= (Me)'(Me) \\ &= y - X(X'X)^{-1}X'y & &= e'Me \\ &= (I_n - X(X'X)^{-1}X')y & &= e'MMe \quad (\text{symmetric } M) \\ &= My & &= e'Me \quad (\text{idempotent } M) \\ &= M(X\beta + e) & &= \text{trace}(e'Me) \quad (e'Me \text{ is a scalar}) \\ &= Me & &= \text{trace}(Me'e') \\ & & &= \text{trace}(Me'e') \end{aligned}$$

$$\text{rank}(M) = \text{trace}(M)$$

$$\begin{aligned} &= \text{trace}(I_n - X(X'X)^{-1}X') \\ &= \text{trace}(I_n) - \text{trace}(X(X'X)^{-1}X') \\ &= n - \text{trace}(X'X(X'X)^{-1}) \\ &= n - \text{trace}(I_k) \\ &= n - k \end{aligned}$$

$$\text{Unbiased estimator of } \sigma^2: S^2 = \frac{1}{n-k} \sum_{i=1}^k \hat{e}_i^2 = \frac{1}{n-k} \hat{e}'\hat{e}.$$

$$\begin{aligned} \text{Then: } V &= \text{Var}(\hat{\beta}_{0:n} | X_1, \dots, X_n) = \sigma^2(X'X)^{-1}. \text{ non stochastic} \\ \hat{V} &= \text{Var}(\hat{\beta}_{0:n}) = S^2(X'X)^{-1} = \hat{\text{Var}}(\hat{\beta}_{0:n} | X). \text{ stochastic.} \end{aligned}$$

Inference on $\hat{\beta}_{0:n}$ (finite sample):

$$\begin{aligned} \text{Assumption: } Y|X=x &\sim N(X\beta, \sigma^2) \text{ for all } x \in \text{supp}(X) \\ \Rightarrow E[Y|X=x] &= X'\beta, \text{Var}(Y|X=x) = \sigma^2, E = Y - X\beta, E \perp X, E \sim N(0, \sigma^2) \end{aligned}$$

$$\text{Then: } y|X \sim N(X\beta, \sigma^2 I_n)$$

$$\text{define } e_i = y_i - x_i'\beta, \text{ then } e_i \sim N(0, \sigma^2), e \sim N(0, \sigma^2 I_n).$$

$$\text{Theorem: } \hat{\beta}_{0:n}|X \sim N(\beta, \sigma^2(X'X)^{-1})$$

$$M = I_n - X(X'X)^{-1}X'$$

$$\left(\begin{array}{c} \hat{\beta}_{0:n} - \beta \\ \hat{e} \end{array} \right) \sim N(0, \begin{pmatrix} \sigma^2(X'X)^{-1} & 0 \\ 0 & \sigma^2 M \end{pmatrix}).$$

$$\text{Corollary: } \hat{\beta}_{0:n} \perp \hat{e} | X, \hat{\beta}_{0:n} \perp S^2 | X.$$

$$\text{Corollary: } \hat{\beta}_k | X \sim N(\beta_k, \sigma^2[(X'X)^{-1}]_{kk})$$

Get CIs for β_k assuming σ^2 is known: marginal

$$H_{\text{null}}: \beta_k = b_k \quad H_{\text{alt}}: \beta_k \neq b_k. \quad (b_k = 0 \text{ most common}).$$

Define the t-statistic:

$$T_k(\beta_k) = \frac{\hat{\beta}_k - b_k}{\text{stderr}(\hat{\beta}_k)} = \frac{\hat{\beta}_k - b_k}{\sqrt{\sigma^2[(X'X)^{-1}]_{kk}}}$$

proposition: Under the null, $T_k(\beta_k) \sim N(0, 1)$ unconditionally.

Conducting α level α hypothesis test:

$$N(0, 1) \text{ symmetric.}$$

① Compute T_k

② Compute $z_{\alpha/2}$ and $z_{1-\alpha/2}$ under $N(0, 1)$. \Leftrightarrow let $C(X) = Z_{1-\alpha/2}$. then:

③ Reject H_{null} if $|T_k| < z_{\alpha/2}$ or $|T_k| > z_{1-\alpha/2}$. $\phi_{t\text{-test}} = \mathbb{1}[|T_k| > C(\alpha)]$

$$E[\hat{e}'\hat{e}|X] = E[\text{trace}(Me'e')|X]$$

$$= \text{trace}(E[Me'e'|X])$$

$= \text{trace}(M E[e'e']|X)$ M is constant given X

$= \text{trace}(M \sigma^2 I_n)$ homoskedasticity & iid

$$= \sigma^2 \text{trace}(M)$$

$$= \sigma^2 \text{rank}(M) \quad M \text{ idempotent}$$

proposition: $E[\hat{V}] = V$.

$$\text{prof: } E[\hat{V}] = E[S^2(X'X)^{-1}]$$

$$= E[E[S^2(X'X)^{-1}|X]]$$

$$= E[E[S^2|X](X'X)^{-1}]$$

$$= E[\sigma^2(X'X)^{-1}]$$

$$= V$$

$$\text{prof: conditional on } X: \left(\begin{array}{c} \hat{\beta}_{0:n} - \beta \\ \hat{e} \end{array} \right) = \begin{bmatrix} (X'X)^{-1}X' \\ M \end{bmatrix} e$$

$$\sim \begin{bmatrix} (X'X)^{-1}X' \\ M \end{bmatrix} N(0, \sigma^2 I_n)$$

remember that $A \cdot N(0, V) = N(0, AVA')$

$$\therefore \text{var}(\begin{bmatrix} \hat{\beta}_{0:n} - \beta \\ \hat{e} \end{bmatrix} | X) = \begin{bmatrix} (X'X)^{-1}X' \\ M \end{bmatrix} \sigma^2 I_n [X(X'X)^{-1} M']$$

$$= \sigma^2 \begin{bmatrix} (X'X)^{-1} & 0 \\ 0 & M \end{bmatrix}$$

Functions of standard normal distributions.

1. χ_n^2 : let $Z \sim N(0, \sigma^2 I_n)$, then

$$Z'Z = \sum_{i=1}^n Z_i^2 \sim \chi_n^2$$

2. F-distribution: let $W_1 \sim \chi_m^2$, $W_2 \sim \chi_n^2$, $W_1 \perp W_2$

$$F = \frac{W_1/m}{W_2/n} \sim F(m, n)$$

3. Student t-distribution: let $Z \sim N(0, 1)$, $W \sim \chi_n^2$, $Z \perp W$

$$U = Z/\sqrt{W/n} \sim t_n \quad U^2 \sim F(1, n)$$

proposition: y : r.n-vector, $y \sim N(0, \Sigma)$, Σ symmetric positive semidefinite full rank. then $y'\Sigma^{-1}y \sim \chi_n^2$

lemma: A, B : r. vectors, if $F_{AB}(a|b) = F_A(a)$ for all $(a, b) \in \mathbb{R}^{\dim(A)} \times \mathbb{R}^{\dim(B)}$ $f_{AB}(a|b) > 0$, then $A \perp B$.

proposition: Exact size control.

Suppose σ^2 known. Let P_0 be the dist. of (Y, X_1, \dots, X_k) such that H_{null} is true. Then $P_0(\phi_{t-\text{test}} = 1) = \alpha$

$$\begin{aligned} \text{proof: } P_0(\phi_{t-\text{test}} = 1) &= P_0(T_k < Z_{\alpha/2} \text{ or } T_k > Z_{1-\alpha/2}) \\ &= P_0(T_k < Z_{\alpha/2}) + P_0(T_k > Z_{1-\alpha/2}) \\ &= \frac{\alpha}{2} + (1 - (1 - \alpha/2)) \\ &= \alpha. \end{aligned}$$

Constructing a $100(1-\alpha)\%$ confidence interval for β_k by inverting the test.

$$T_k(b_k) = \frac{\hat{\beta}_k - b_k}{\text{stderr}(\hat{\beta}_k)}$$

$$\begin{aligned} I_{\text{in}(1-\alpha)} &= \{b_k \in \mathbb{R} : Z_{1-\alpha/2} \leq T_k(b_k) \leq Z_{\alpha/2}\} \\ &= [\hat{\beta}_k - Z_{1-\alpha/2} \text{ stderr}(\hat{\beta}_k), \hat{\beta}_k + Z_{\alpha/2} \text{ stderr}(\hat{\beta}_k)] \\ &= [\hat{\beta}_k - C(\alpha) \text{ stderr}(\hat{\beta}_k), \hat{\beta}_k + C(\alpha) \text{ stderr}(\hat{\beta}_k)] \end{aligned}$$

$$P(I_{\text{in}(1-\alpha)} \ni \beta_k) = 1 - \alpha$$

$$\begin{aligned} |T_k| < C(\alpha) &\quad -C(\alpha) \leq \frac{\hat{\beta}_k - b_k}{\text{stderr}(\hat{\beta}_k)} \leq C(\alpha) \\ &\quad -C(\alpha) \leq \frac{b_k - \hat{\beta}_k}{\text{stderr}(\hat{\beta}_k)} \leq C(\alpha) \\ \hat{\beta}_k - C(\alpha) \text{ stderr}(\hat{\beta}_k) \leq b_k &\leq \hat{\beta}_k + C(\alpha) \text{ stderr}(\hat{\beta}_k) \\ \hat{\beta}_k - Z_{1-\alpha/2} \text{ stderr}(\hat{\beta}_k) \leq b_k &\leq \hat{\beta}_k + Z_{\alpha/2} \text{ stderr}(\hat{\beta}_k) \end{aligned}$$

Get CIs for β assuming σ^2 is known. Joint.

Simple idea: suppose $k=2 \Rightarrow \beta_1, \beta_2$; then: $C\text{I}^{\text{joint}} = (I_1(1-\alpha) \times I_2(1-\alpha))$ get a box.

suppose that $\{I_1(1-\alpha) \ni \beta_1\}$ and $\{I_2(1-\alpha) \ni \beta_2\}$ are independent. then: $P(C\text{I}_1 \times C\text{I}_2 \ni (\beta_1, \beta_2)) = P(C\text{I}_1 \ni \beta_1) \cdot P(C\text{I}_2 \ni \beta_2)$

Marginal \rightarrow Joint WRONG!

$$\begin{aligned} &= (1-\alpha)^2 \\ &\neq (1-\alpha) \end{aligned}$$

The Wald statistic based confidence set.

$$\hat{\beta} - \beta | X \sim N(0, \sigma^2(X'X)^{-1}) = \sigma(X'X)^{-1} N(0, I_n)$$

$$\therefore \sigma(X'X)^{-1} (\hat{\beta} - \beta) | X \sim N(0, I_n)$$

$$\therefore \text{define } W(\beta) \equiv (\hat{\beta} - \beta)' \sigma^2(X'X)^{-1} (\hat{\beta} - \beta) \sim \chi_k^2$$

define $C(\alpha)$ to be $1-\alpha$ quantile of χ_k^2 . $G(C(\alpha)) = 1-\alpha$, $G(\cdot)$ is the cdf of χ_k^2 .

$$\text{define } C\text{I}_n^W(1-\alpha) = \{b \in \mathbb{R}^k : W(b) \leq C(\alpha)\}$$

proposition: Exact coverage: $P(C\text{I}_n^W(1-\alpha) \ni \beta) = 1 - \alpha$ shape: confidence ellipsoids.

$$\begin{aligned} \text{proof: } P(C\text{I}_n^W(1-\alpha) \ni \beta) &= P(W(\beta) \leq C(\alpha)) \\ &= G(C(\alpha)) \\ &= 1 - \alpha. \end{aligned}$$

Inference with unknown σ^2

H_{null} : $\beta_k = b_k$ H_{alt} : $\beta_k \neq b_k$

$$T_k = \frac{\hat{\beta}_k - b_k}{\text{stderr}(\hat{\beta}_k)} = \frac{\hat{\beta}_k - b_k}{\sqrt{s^2(X'X)^{-1}}_{kk}}$$

Theorem: Under the null $\beta_k = b_k$, $T_k \sim t_{n-k}$

$$W(\beta) \equiv (\hat{\beta} - \beta)' s^{-2}(X'X)(\hat{\beta} - \beta)$$

$$W(\beta) \sim \frac{\chi_k^2}{\chi_{n-k}^2 / (n-k)} \quad F(\beta) = \frac{W(\beta)}{k} \sim \frac{\chi_k^2 / k}{\chi_{n-k}^2 / (n-k)}$$

proposition: Under the null $\beta = b$: $F(b) \sim F(k, n-k)$

WLS

$$\hat{\beta}_{\text{WLS}}(w) = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n w_i (y_i - x_i' b)^2 \quad w_i = w(x_i)$$

$$w = \begin{bmatrix} w_1 & \cdots & 0 \\ 0 & \cdots & w_n \end{bmatrix}$$

$$\hat{\beta}_{\text{OLS}}(w) = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} (y - Xb)' w (y - Xb)$$

$$\begin{aligned} \text{proof: objective function: } (y - Xb)' w (y - Xb) &= (y - Xb)' w (y - Xb) \\ &= y' w y - y' w X' b - b' X' w y - b' X' w X b \\ \text{take derivative wrt } b: \quad -2X' w y + 2X' w X b &= 0 \\ \therefore \hat{\beta}_{\text{OLS}} &= (X' w X)^{-1} X' w y \end{aligned}$$

gives: $\hat{\beta}_{WLS}(w) = (X'wX)^{-1}X'wY$

Theorem: $\hat{\beta}_{WLS} \xrightarrow{P} E[Y|X=x]^T E[X] \equiv \hat{\beta}$ It is a generalization of OLS estimator

Theorem: If $E[Y|X=x] = X'\beta$, then $\hat{\beta} = \beta$

Finite Sample Properties of WLS.

Proposition unbiased under linearity & variance

Suppose ① linearity: $E[Y|X=x] = X'\beta$

② iid $\{y_i, x_i\} \sim \text{③ } w_i = w(x_i)$

then: $E[\hat{\beta}_{WLS}] = \beta$

$$\text{Var}(\hat{\beta}_{WLS}|X) = (X'wX)^{-1}X'w\Omega wX(X'wX)^{-1}$$

$$\Omega = \begin{bmatrix} \sigma^2(x_1) & & \\ & \ddots & 0 \\ 0 & & \sigma^2(x_n) \end{bmatrix} \text{ where } \sigma^2(x) = \text{Var}(Y|X=x).$$

If we pick $w = \Omega^{-1}$: then $\text{Var}(\hat{\beta}_{WLS}|X) = (X'\Omega^{-1}X)^{-1}X'\Omega^{-1}\Omega\Omega^{-1}X(X'\Omega^{-1}X)^{-1} = (X'\Omega^{-1}X)^{-1}$

$$w = \Omega^{-1} = \begin{bmatrix} 1/\sigma^2(x_1) & & 0 \\ 0 & \ddots & 1/\sigma^2(x_n) \end{bmatrix}$$

(infeasible) GLS: $\hat{\beta}_{GLS} = \hat{\beta}_{WLS}(\Omega^{-1})$

* $\text{Var}(\hat{\beta}_{GLS}) \leq \text{Var}(\hat{\beta})$ & $\hat{\beta}$ that is linear and unbiased conditional on X .

$$\hat{\beta}_{GLS} = \hat{\beta}_{WLS}(w_i = \frac{1}{\sigma^2(x_i)}) = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma^2(x_i)} (y_i - x_i'b)^2$$

feasible GLS: FGLS:

$$\sigma^2(x) = \text{Var}(E[Y|X=x]) = E[E^2|X=x]$$

$$\text{Assume } \sigma^2(x) = \exp(\alpha'x)$$

$$\text{① get } \hat{e}_i = y_i - x_i'\beta_{OLS}$$

$$\text{② get } \hat{\alpha} \text{ from WLS of } \hat{e}^2 \text{ on } x: \hat{\alpha}_{FGLS} = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (\hat{e}_i^2 - \exp(\alpha'x_i))^2$$

$$\text{③ compute } \hat{\sigma}^2(x_i) = \exp(\hat{\alpha}'x_i)$$

Problem: ① requires estimation of $\sigma^2(x)$, functional form not correct.

Asymptotic Properties of OLS Estimator

Theorem: Consistency.

① iid sample ② finite moments ③ sufficient variation.

then: $\hat{\beta}_{OLS} \xrightarrow{P} \beta_{OLS}$

Theorem: Asymptotic Normal

$$\text{define } D = E[XX'] (= \text{Var}(x)) = E[E^2 XX']$$

$$\text{then } \sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{d} N(0, D^{-1}CD^{-1})$$

Estimating the asymptotic covariance matrix:

$$\sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{d} N(0, D^{-1}CD^{-1})$$

$$V \equiv D^{-1}CD^{-1} = E[XX']^{-1}E[E^2 XX']E[XX']^{-1}$$

$$\hat{V} \equiv \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{e}_i^2 x_i x_i' \right) \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1}$$

\hat{V} heteroskedasticity robust covariance matrix
(std. err.) estimator.

$$\text{proof: } \hat{\beta}_{OLS} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) \text{ by OLS} \xrightarrow{P} E[XX']^{-1} E[XY]$$

$$\begin{aligned} \text{proof: } \hat{\beta}_{OLS} &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) \\ &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i (x_i'\beta_{OLS} + e_i) \right) \\ &= \beta_{OLS} + \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i e_i \right) \end{aligned}$$

$$\therefore \hat{\beta}_{OLS} - \beta_{OLS} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i e_i \right)$$

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{OLS} - \beta_{OLS}) &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left[\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n x_i e_i \right) \right] \\ &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \underbrace{\left[\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n x_i e_i \right) - E[x_i e_i] \right]}_{\textcircled{1}} \end{aligned}$$

$$\text{①} \xrightarrow{d} N(0, \text{Var}(x)e) = N(0, C)$$

$$\therefore \sqrt{n}(\hat{\beta}_{OLS} - \beta_{OLS}) \xrightarrow{d} D^{-1}N(0, C) \text{ by CLT.}$$

$$\sqrt{n}(\hat{\beta}_{OLS} - \beta_{OLS}) \xrightarrow{d} N(0, D^{-1}CD^{-1})$$

Why this result is important?

$$\text{Remember. } \text{var}(\hat{\beta}_{\text{out}} | X) = (X'X)^{-1} X' \Omega X (X'X)^{-1}$$

$$\Omega = \begin{bmatrix} \sigma^2(x_1) & & 0 \\ & \ddots & \\ 0 & & \sigma^2(x_n) \end{bmatrix}$$

$$\begin{aligned} n \text{var}(\hat{\beta}_{\text{out}} | X) &= n(X'X)^{-1} X' \Omega X (X'X)^{-1} \\ &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i^2 x_i x_i' \right) \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \end{aligned}$$