

Sample Space & Events

- Sample Space S is a collection of possible outcomes of a trial
- An event A is a subset of all possible outcomes

* Events are sets

$$\boxed{1.} \subset = \text{subset}$$

$$\boxed{2.} \cap = \text{Intersection}$$

$$\boxed{3.} \cup = \text{Union}$$

Properties:

$$\bullet A \subset B, A \cup B = B$$

$$\bullet A \subset B \& B \subset A, A = B$$

$$\bullet A \cup A^c = S$$

↳ complement set

Probability

- Every event A assigned Probability we require:

$$\boxed{1.} P(A) \geq 0 \text{ for all } A \subset S$$

$$\boxed{2.} P(S) = 1$$

$$\boxed{3.} \text{For disjoint sets } A_1, A_2, P(A_1 \cup A_2) = \sum P(A_i)$$

Properties of probability:

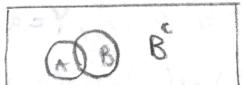
$$\bullet P(A^c) = 1 - P(A)$$

$$\bullet P(\emptyset) = 0$$

$$\bullet \text{if } A \subset B, P(A) \leq P(B)$$

$$\bullet P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$\bullet P(A \cap B^c) = P(A) - P(A \cap B)$$



Simple Sample Space

- All events Equally Likely

$$P(A) = n(A) / n(S)$$

Combinatorics

- If experiment happens in multiple parts

$$\left. \begin{array}{l} \text{part 1: } N \text{ outcomes} \\ \text{part 2: } M \text{ outcomes} \end{array} \right\} \text{Total: } M \times N$$

- Permutation → # of Ordered ways to count objects

$$\text{K objects} \left\{ \frac{N!}{(N-K)!}$$

14.30 Summary Sheet

- Combination → # of unordered ways to count objects

$$K \text{ objects out of } N \left\{ \frac{N!}{(N-K)! K!} = \binom{N}{K} \right.$$

Independence

- Events A & B independent

$$\text{if } P(AB) = P(A)P(B)$$

↳ Knowing one event occurs doesn't give you intuition about the other

$$\begin{aligned} \text{Probability } A \text{ happens} \\ \text{exactly } k \text{ times in next } N \text{ trials} \end{aligned} \left\{ \frac{N!}{(N-k)! k!} P(A)^k P(A^c)^{N-k} \right.$$

Conditional Probability

- One event happening does tell you something about the other

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Versions of Bayes' Theorem

$$\boxed{1.} P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$\boxed{2.} P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)}$$

Random Variables

- X a random variable is a real-valued function whose domain is the sample space

$$\text{Binomial Distribution} \left\{ f_x(x) = \binom{n}{x} p^x (1-p)^{n-x} \right.$$

Probability Function (Discrete RV)

- Discrete random Variable X

$$\text{Probability Function} \left\{ f_x(x) = P(X=x) \right.$$

Properties:

$$\begin{aligned} \boxed{1.} 0 \leq f_x(x) \leq 1 \\ \boxed{2.} \sum_x f_x(x) = 1 \end{aligned} \quad \boxed{3.} P(A) = \sum_{x \in A} f_x(x)$$

Probability Density Function (PDF)

- Function for Continuous RV that x lies in particular region

$$P(X \in A) = \int_A f_x(x) dx$$

Cumulative Distribution Function (CDF)

- CDF of Random Variable x is

$$F_x(x) = P(X \leq x)$$

Properties:

$$\boxed{1.} F_x(x) \text{ is non-decreasing in } x$$

$$\boxed{2.} \lim_{x \rightarrow -\infty} F_x(x) = 0 \quad \boxed{3.} \lim_{x \rightarrow \infty} F_x(x) = 1 \quad \left\{ \begin{array}{l} F_x(x) = \int_{-\infty}^x f_x(t) dt \\ F'(x) = f_x(x) \end{array} \right.$$

CDF → Discrete RV

- Always Right-continuous

$$\text{ex. } F_x(x) = \begin{cases} 0, & x < 0 \\ 1-p, & 0 \leq x \leq 1 \\ 1, & x \geq 1 \end{cases}$$

example:

$$f_x(x) = \begin{cases} 1, & \{a < x < b\} \\ 0, & \text{else} \end{cases} \frac{1}{b-a}$$

↳ Indicator Function

$$F_x(x) = \begin{cases} 0, & x < a \\ (x-a)/(b-a), & a \leq x \leq b \\ 1, & x > b \end{cases}$$

Normal Distribution

$$f_x(x) = \phi(x) = \frac{1}{\sqrt{2\pi}} * e^{-x^2/2}$$

Joint Distributions

$$\text{Function: } P((x,y) \in A) = \iint f_{xy}(x,y) dx dy$$

Joint + Marginal Distributions

- we can Recover Marginal distributions from Joint

$$\text{Discrete} \left\{ f_x(x) = \sum_y f_{xy}(x,y) \right.$$

$$\text{Continuous} \left\{ f_x(x) = \int_y f_{xy}(x,y) dy \right.$$

Independence of RW

- If \boxed{X} & \boxed{Y} are independent:

$$F_{xy}(x,y) = F_x(x)F_y(y)$$

If continuous:

$$f_{xy}(x,y) = f_x(x)f_y(y)$$

Conditional Distribution

- Conditional PDF of \boxed{Y} given \boxed{X}

$$f_{y|x}(y|x) = f_{xy}(x,y)/f_x(x)$$

(when discrete)

$$f_{y|x}(y|x) = P(Y=y|X=x)$$

Example

$$f_{xy}(x,y) = \begin{cases} 2, & 0 \leq y \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$f_{x|y}(x|y) = \frac{1}{1-y} \{y \leq x \leq 1\}$$

$$f_{y|x}(y|x) = \frac{1}{x} \{0 \leq y \leq x\}$$

conditional Distribution / Independence

- If independent

$$f_{y|x}(y|x) = f_y(y) \quad \boxed{8} \quad f_{xy} = f_x f_y$$

Functions of Random Variables

- \boxed{X} is random variable where $f_x(x)$ known.

we want $Y = h(x)$

$$F_y(y) = \int_{\{x : h(x) \leq y\}} f_x(x) dx$$

↳ to find CDF integrate appropriate region

$$\text{If } Y \text{ continuous: } f_y(y) = F'_y(y)$$

Example: Function of RV

$$f_x(x) = \begin{cases} 1/2, & \text{for } -1 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$\boxed{Y = X^2}$, what is $f_y(y)$?

$$\text{steel: } F_y(y) = P(Y \leq y) \\ = P(X^2 \leq y) \\ = P(-\sqrt{y} \leq X \leq \sqrt{y})$$

$$\text{so } F_y(y) = \int_{-\sqrt{y}}^{\sqrt{y}} \frac{1}{2} dx = \left[-\sqrt{y} \right]$$

$$F_y(y) = \begin{cases} 0, & y < 0 \\ \sqrt{y}, & 0 \leq y \leq 1 \\ 1, & y > 1 \end{cases}$$

$$f_y(y) = \begin{cases} \frac{1}{2\sqrt{y}}, & 0 \leq y \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

Linear Transformation Example

- \boxed{X} has PDF $f_x(x)$

↳ Let $\boxed{Y = a + bX}$ (How distribute)

case 1: $\boxed{b > 0}$

$$F_y(y) = P(Y \leq y) = P(a+bX \leq y) \\ = P(X \leq \frac{(y-a)}{b}) \\ = \int_{-\infty}^{\frac{(y-a)}{b}} f_x(x) dx$$

$$\text{so } \boxed{\text{PDF} = f_y(y) = F'_y(y) = \frac{1}{b} f_x\left(\frac{y-a}{b}\right)}$$

Expectation (or Mean)

continuous: $E(X) = \int x f_x(x) dx$

$$\text{discrete: } E(X) = \sum_{j=1}^J x_j f_x(x_j)$$

* Properties of Expectation

$$1. E(a) = a, \text{ for constant}$$

$$2. E(Y) = a + bE(X), Y = a + bX$$

$$3. E(Y) = E(X_1) + E(X_2) + \dots + E(X_n)$$

↳ X doesn't need independence

$$4. E(XY) = E(X)E(Y) \text{ if } \text{Independent}$$

Expectation of RV Function

$$\text{we know } \boxed{f_x(x)}, Y = g(x)$$

we could figure out how Y is distributed & find $f_y(y)$

$$\text{so } \boxed{E(Y) = \int y f_y(y)}$$

Easier,

$$E(Y) = E(g(x)) = \int g(x) f_x(x)$$

Variance

$$1. \text{Var}(x) = E[(x - \mu)^2]$$

$$\text{so } \text{Var}(x) = \sigma_x^2$$

Properties of Variance

$$1. \text{Var}(x) \geq 0$$

$$2. \text{Var}(a) = 0$$

$$3. \text{Var}(Y) = b^2 \text{Var}(x), Y = a + bX$$

$$4. \text{Var}(Y) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)$$

↳ X_1, X_2, \dots, X_n independent

$$5. \text{Var}(x) = E[x^2] - (E[x])^2$$

**

Variance of a Function

WANT: $\text{Var}(Y)$, where $Y = r(x)$

$$\text{Var}(Y) = E[Y^2] - E[Y]^2$$

$$= \int r(x)^2 f_x(x) dx - \left[\int r(x) f_x(x) dx \right]^2$$

Exam 1

X takes on 0 or 1 (Bernoulli RV)

$$P(X=0) = 0.22 \quad P(X=1) = 0.78$$

\bar{X}_n depends on n

if $n=1$ if $n=2$

$$\bar{X}_1 = 0, 1 \quad P(\bar{X}_1 = 0) = 0.22^2 = 0.0484$$

$$P(\bar{X}_1 = 1) = 0.78 \quad P(\bar{X}_1 = 0.5) = 34.32$$

$$P(\bar{X}_1 = 1) = 0.684$$

Central Limit Theorem

Let X_1, \dots, X_n be a random sample (i.i.d) of size n w/ finite mean μ & variance σ^2

$$\lim_{n \rightarrow \infty} P\left[\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq x\right] = \Phi(x)$$

$\hookrightarrow \Phi(x)$ = Normal CDF

Estimator

What makes "good" Estimator

$\hookrightarrow \hat{\theta}$ is unbiased if & only if

$$E[\hat{\theta}] = \theta_0$$

$\hookrightarrow \hat{\theta}$ is consistent if for all $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta_0| > \epsilon) = 0$$

(Law of Large Numbers)

* Sample mean is consistent estimator of population mean

More Criteria

If estimator not consistent, it's bad

\hookrightarrow Let $n/2$ be largest integer less than $n/2$

$\tilde{\theta} = \frac{1}{[n/2]}$, this is unbiased & consistent, why prefer \bar{X} ?

$$\text{Var}(\tilde{\theta}) = \frac{2\sigma^2}{n} \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

\hookrightarrow we say $\hat{\theta}$ is more efficient than $\tilde{\theta}$

Efficiency of Estimators

efficiency comparison hard to do for most

Many Estimators are consistent w/ Normal Distribution in large samples

Compare them on Variance of Large Sample Standard Error

Confidence Intervals

• corresponds to margin of error in polls

say that, $\hat{p} = 0.42$ & $n = 1500$

standard

Error

$$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

a 95% confidence interval is

$$CI_{95} = [\hat{p} - 1.96 \cdot SE(\hat{p}), \hat{p} + 1.96 \cdot SE(\hat{p})]$$

$$\hookrightarrow P(\hat{p} \in CI_{95}) \rightarrow 0.95 \text{ as } n \rightarrow \infty$$

Properties of Confidence Interval

• It's a random interval

\hookrightarrow contains true value of parameter with some probability

• Data changes \rightarrow But interval stays the same

Z-score

$$Z_n = \frac{\hat{p} - E[\hat{p}]}{SD(\hat{p})} \quad \begin{cases} \hat{p} = p_0 \\ SD(\hat{p}) = \sqrt{\frac{p_0(1-p_0)}{n}} \end{cases}$$

2 Properties:

1. \hat{p} is sample mean so by CLT

$$\hookrightarrow Z_n \sim N(0,1)$$

2. $P_0 \in CI_{95}$ if $-1.96 \leq Z_n \leq 1.96$

For General Distributions

$$SD(\hat{x}) = \sigma_x / \sqrt{n}$$

$$\frac{\bar{X} - \mu_x}{SD(\bar{X})} = \frac{\bar{X} - \mu_x}{\sigma_x / \sqrt{n}} \sim N(0,1)$$

$$\hat{\sigma}_x^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

• If c_a is the critical value for $N(0,1)$ satisfying

$$P(-c_a < N(0,1) < c_a) = 0.95$$

$$CI_d^{H_x} = [\bar{X} - c_a \sqrt{\frac{s^2}{n}}, \bar{X} + c_a \sqrt{\frac{s^2}{n}}]$$

Hypothesis Tests

• consider Parameter $\mu = E[Y]$

we test a Null Hypothesis about μ

$$H_0: \mu = m$$

against alternative hypothesis $H_1: \mu \neq m$

• we (reject or don't reject) the Null

Discrete Estimators

$$\text{sample Average} \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\text{standard Error} \quad \text{s.e.}(\bar{Y}) = \sqrt{\frac{\sigma^2}{n}}$$

$$\hookrightarrow \text{Var} = \sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

+ Statistic

$$t = \frac{\bar{Y} - m}{\text{s.e.}(\bar{Y})} \stackrel{d}{\rightarrow} N(0,1)$$

$|t| > 1.96$, Reject the null

$|t| < 1.96$, Don't reject Null

Equivalently

$$p = Pr(|t| > 1.96)$$

$p < 0.05$ reject Null \Rightarrow you
 $p \geq 0.05$ Don't reject Null

Terminology

test statistic \rightarrow stat used for testing null (t)

Null Distribution \rightarrow sample distribution

Type I error \rightarrow rejecting Null even if true

Type II error \rightarrow Not rejecting Null if false

Critical Region \rightarrow set of values such that $|t| \in \text{Critical Region}$ we reject null (1.96, 0)

Two Sample Test

Two Samples (X_1, \dots, X_n) & $(X_{n+1}, \dots, X_{n+m})$

$$\begin{cases} H_0: \mu_1 = \mu_2 = \bar{\mu} \\ H_1: \mu_1 \neq \mu_2 \end{cases} \quad \begin{cases} \text{test to see if same mean} \end{cases}$$

Test statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}}$$

Conditional Expectation

- Is the expectation of a conditional distribution of Y given X

$$E(Y|X=x) = \int y f_{Y|X}(y|x) dy$$

Law of Iterated expectations

$$E[Y] = E[E[Y|X]]$$

Conditional Variance

- $\text{Var}(Y|X=x)$ is Variance of Distribution of Y given X

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

Law of Total Variance:

$$\text{Var}(Y) = E(\text{Var}[Y|X]) + \text{Var}(E[Y|X])$$

Conditional Variance: Example

- N = # of units produced in a year
- $E(N) = 2$, $\text{Var}(N) = 2$

$\hookrightarrow p_{\text{success}} = 0.2$

Probability of 5 successes

$$\binom{5}{3} (0.2)^3 (1-0.2)^2 = 0.0512$$

of expected successes

$$E(S) = E(E(S|N)) = E(Np) = 0.2(E(N)) = 0.4$$

Covariance & Correlation

- Key Moment of Joint Distribution is Covariance

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}[X, Y]}{\sigma_x \sigma_y}$$

\Downarrow

$$\rho_{xy}$$

$P > 0$	$X \& Y$	Positively correlated
$P < 0$	$X \& Y$	Negatively correlated
$P = 0$	Uncorrelated	

Properties of Covariance & P

- $\text{Cov}(X, X) = \text{Var}(X)$
- $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$
- $\text{Cov}(X, Y) = 0$
 \hookrightarrow if $X \& Y$ independent
- $\text{Cov}(aX+b, cY+d) = ac(\text{Cov}(X, Y))$
- $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$

- $| \text{Cov}(X, Y) | \leq 1$
- $| \text{Cov}(X, Y) | = 1$, if $Y = a + bX$

Other Moment Metrics

mean = $E(Y)$

Variance = $E[(Y - \mu_Y)^2]$

skewness = $E\left[\frac{(Y - \mu_Y)^3}{\sigma^3}\right]$

Kurtosis = $E\left[\frac{(Y - \mu_Y)^4}{\sigma^4}\right]$

Mode = Point where PDF has highest value

Median = Point where Integral of PDF = 0.5

Median

Measure of Location where

$$F_X(m) = \frac{1}{2} \quad \begin{matrix} * \text{halfway} \\ \text{through dataset} \end{matrix}$$

Quantiles

- Let $a \in (0, 1)$, CDF strictly increasing

- a^{th} quantile of X is $\alpha_a(x)$ such that

$$F(\alpha_a(x)) = a$$

- 4 quartiles
- 5 quintiles $\sim 0.9^{\text{th}}$ quantile = 90th percentile

$$\alpha_x(0.75) - \alpha_y(0.5) = \text{interquantile range}$$

\hookrightarrow measure of Dispersion

Properties of Median & Quantiles

$$\text{med}(bX+c) = b * \text{med}(x) + c$$

$$\text{med}(X+Y) \neq \text{med}(x) + \text{med}(y)$$

In General:

$$\alpha_{bX+c} = b \alpha_x + c$$

$$\alpha_{X+Y} \neq \alpha_x + \alpha_y$$

Markov Inequality

- X is RV that non-negative

$\hookrightarrow E[X]$ exists

For any $t > 0$,

$$P(X \geq t) \leq E(X)/t$$

Chebychev Inequality

- X is a random variable for which $\text{Var}(x)$ exists

For any $t > 0$,

$$P(|X - E(x)| \geq t) \leq \text{Var}(x)/t^2$$

Simple Random Sampling

- Choose an individual at Random from population

1. X_1, X_2 are independent

2. X_1, X_2 are identically distributed

Dataset $\{X_1, \dots, X_n\}$ are i.i.d.

Statistics

- Assume we have data set $\{x_1, x_n\}$

\hookrightarrow each x_i has $F_x(x)$

- A Parameter characterises $F_x(x)$

1. Expectation $E(x)$

2. Variance $\text{Var}(x)$

3. $\theta_0 = \text{True Parameter}$

\hookrightarrow Reflects F_x data drawn from

* goal is to learn θ_0 from data & quantify uncertainty

Estimator

\hookrightarrow uses data to construct a "guess" of parameter value

$$\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$$

\hookrightarrow function of data

Sample Mean

- Arithmetic average of n random variables

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

Linear Algebra Review

- A Matrix is a rectangular array

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1K} \\ a_{21} & a_{22} & \dots & a_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mK} \end{pmatrix}$$

- Only Matrices of same dimensions can be added/subtracted

- A scalar c can be multiplied w/ matrix

- Transpose of a matrix achieved by turning Rows \rightarrow Columns

$A = \text{Matrix}$

$A' = \text{Transpose of Matrix}$

- a Column vector \rightarrow Matrix of single column

- Matrices must be conformable to be multiplied

- Transpose of a Product

$$(AB)' = B'A'$$

Square Matrices

- Square Matrix has same # of rows & columns

- Square matrix is symmetric if for all elements $a_{ij} = a_{ji}$

- Diagonal Matrix is square matrix w/ only zeroes outside Main Diagonal

Properties of Matrices

- Square Matrix is singular if all its columns are linearly independent

- If non-singular A^{-1} exists system of linear equations

$$Ax = b \Rightarrow x = A^{-1}b$$

Random Vectors

- Let X be a $(k \times 1)$ Random Vector

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_K \end{pmatrix} \quad \text{vector of Means} \quad M = \begin{pmatrix} M_1 \\ M_2 \\ \vdots \\ M_K \end{pmatrix}$$

- Variance of X is a $(K \times K)$ matrix

14.30 Summary Sheet #2

$$\text{Var}(X) = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_K) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) & & \\ \vdots & & \ddots & \\ \text{Cov}(X_K, X_1) & & \text{Var}(X_K) & \end{pmatrix}$$

$$\text{Var}(X) = E[(X - \mu)(X - \mu)']$$

B Variance-Covariance Matrix of X

Properties of Random vectors

- $A = (m \times k)$ matrix
- $X = (m \times 1)$ random vector

$$E[AX] = A E[X], \text{Var}(AX) = A \text{Var}(X) A'$$

Example: Functions of Random Vectors

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, E[X] = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \text{Var}(X) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\text{Now Let: } Y_1 = X_1 - X_2 = \begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = AX,$$

$$A = \begin{pmatrix} 1 & -1 \end{pmatrix}$$

$$E[Y] = E[AX] = \begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\text{Var}[Y] = \begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

Normal Random Vectors

- Let X be a $(k \times 1)$ random vector

$$X \sim N(\mu, V)$$

- X is vector of jointly

$\mu = \text{Mean}$

$V = \text{Var-Cov Matrix}$

The Chi-Square (χ^2) Distribution

- Let X_1, \dots, X_K be independent and $N(0, 1)$. Then,

$$Z_K = X_1^2 + \dots + X_K^2 \sim \chi_K^2$$

$$E[Z_K] = K \quad \text{Var}(Z_K) = 2K$$

Normal Random Vector Properties

- If (X_1, X_2) are jointly normal & uncorrelated, they are independent

- If $X \sim N(\mu, V) \Rightarrow Y = AX$

$$Y \sim N(A\mu, AVA')$$

$$\boxed{3.} \text{ If } X \sim N(0, V) \text{ then } Z = X'V^{-1}X$$

Convergence of Random Vectors

- Let U_1, \dots, U_K be $(k \times 1)$ vectors drawn independently such that $E[U_i] = 0$ then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n U_i \xrightarrow{d} N(0, V)$$

$$\hookrightarrow \text{where } V = E[UU']$$

- if \hat{A} is a conformable matrix

$$\hat{A} \frac{1}{\sqrt{n}} \sum_{i=1}^n U_i \xrightarrow{d} N(0, AVA')$$

Population Regression

- Population Regression is the conditional expectation of Y given X

- when Evaluated at particular X

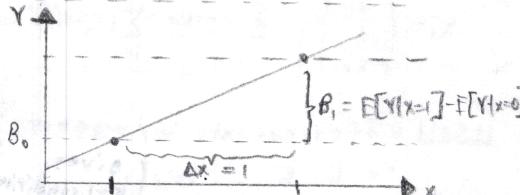
$$E[Y|X=x]$$

Law of Iterated Expectations

$$E(Y) = E(E(Y|X))$$

- If X is Binary, population regression is linear

$$E[Y|X] = \beta_0 + \beta_1 X$$



For Any X :

$$\beta_1 = E[Y|x+1] - E[Y|x]$$

- Similarly, For $E[Y|x] = \beta_0 + \beta_1 X_1 + \beta_2 X_2$

$$\beta_1 = E[Y|x_1=x_1+1, x_2=x_2] - E[Y|x_1=x_1, x_2=x_2]$$

Population Regression Function

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_K \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_K \end{pmatrix}$$

$$\underline{E[Y|X] = X^T \beta}$$

Definition of Error Term

$$H = Y - X^T \beta$$

$$Y = X^T \beta + u \quad E[H|X] = 0$$

$\hookrightarrow Y$ is outcome

$\hookrightarrow X$'s are explanatory variables
aka regressors

β as solution of Minimization Problem

Let,

$$m = E[Y] \quad \& \quad m_x = E[Y|X]$$

Property of Expectations

$$E[(Y-m)^2] \leq E[(Y-c)^2]$$

so if $E[Y|X] = X^T \beta$ then

$$\star \quad \star \quad \star \quad \boxed{\beta = \arg \min_{\beta \in \mathbb{R}^K} E[(Y - X^T \beta)^2]}$$

$\boxed{\beta}$ is argument that minimizes $E[(Y - X^T \beta)^2]$

Estimation: Least Squares

Suppose there exists Random Sample

$$(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)$$

$$\boxed{8} \quad X_i = \begin{pmatrix} X_{i1} \\ \vdots \\ X_{iK} \end{pmatrix} \quad \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^K} \sum_{i=1}^n (Y_i - X_i^T \beta)^2$$

step 1 Differentiate w/ respect to $\beta_0, \beta_1, \dots, \beta_K$ (gives K+1 conditions)

$$\sum_{i=1}^n X_i(Y_i - X_i^T \hat{\beta}) = 0$$

\Downarrow

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X_i^T \right)^{-1} \sum_{i=1}^n X_i Y_i$$

• Co-efficients of $\hat{\beta}$ are called OLS coefficients

Bivariate Regression

$$E[Y|X] = \beta_0 + \beta_1 X$$

OLS Coefficients

$$\hat{\beta}_1 = \frac{\text{cov}(Y_i; X_i)}{\text{var}(X_i)} \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\text{Residuals: } \hat{u}_i = Y_i - \hat{Y}_i$$

Least Squares Simplified Notation

$$y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & X_{11} & \dots & X_{1K} \\ 1 & X_{21} & \dots & X_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \dots & X_{nK} \end{pmatrix} \quad u = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$$

$$y = X \beta + u \quad E[u|X] = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Properties of Least Squares

• Regression Residuals \hat{u}_i have sample mean = 0

$$\frac{1}{n} \sum_{i=1}^n \hat{u}_i = 0$$

• Regressors (X_i) & regression residuals are perfectly un-coordinated

Multi-Collinearity

• Refers to case where columns of X are linearly dependent

in this case $X^T (y - X \beta) = 0$ has many solutions

Unbiased $\hat{\beta}$ for β

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ &= (X^T X)^{-1} X^T (X \beta + u) \\ &= \beta + (X^T X)^{-1} X^T u \quad \downarrow \quad \hat{\beta} \text{ unbiased for } \beta \\ E[\hat{\beta}] &= \beta \end{aligned}$$

Example: log & sin corr

Let $g_{dec} = \log GDP$ per cap gini = log GiniIndex

$$\text{suppose: } E[gini|gdec] = \beta_0 + \beta_1 gdec$$

$$Y_i = \text{gini}_i \quad X_i = \begin{pmatrix} 1 \\ gdec_i \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X_i^T \right)^{-1} \sum_{i=1}^n X_i Y_i \Rightarrow \hat{\beta}_1 =$$

R^2 Metric

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

Large Sample Distribution

Recall

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X_i^T \right)^{-1} \sum_{i=1}^n X_i Y_i$$

$$\sqrt{n}(\hat{\beta} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n X_i X_i^T \right)^{-1} * \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i (Y_i - \beta)$$

(LLN) \downarrow CLT \downarrow

$$\hat{Q} = E[X_i X_i^T]^{-1} \quad N(0, V)$$

$$\therefore \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, V) \quad \hookrightarrow \text{w/ } V = Q^{-1} \sum Q$$

LSD: Estimator Form

$$\hat{\beta} \sim N(0, V/n)$$

• we estimate V as

$$\hat{V} = \hat{Q}^{-1} \sum \hat{Q}$$

$$\hat{Q} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i (Y_i - \hat{\beta})^T X_i$$

• now we can test hypothesis about β

Tests

$$\left| \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} \right| \geq 1.96 \quad \text{where } \text{se}(\hat{\beta}_1) = \sqrt{\hat{V}_{11}/n}$$

homoskedasticity vs. heteroskedasticity

Homo: $\text{Var}(Y|X) = \sigma^2$

Hetero: $\text{Var}(Y|X) = \sigma^2(x)$ depends on X

• homoskedasticity

$$\left| \Sigma = \sigma^2 Q \right| \quad V = \sigma^2 Q^{-1}$$

Interpretation of Regression Coefficients

① $Y = \alpha_0 + \alpha_1 X_1 + V \quad E[V|X_1] = 0$

② $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + H \quad E[H|X_1, X_2] = 0$

α_1 : change in Y w/ unit X_1

β_1 : change in Y w/ unit X_1 , keeping X_2 constant

Omitted Variable Bias

• want to know β_1 but have no data on X_2

\hookrightarrow can only estimate β_1 . What is β_1 ?

$$\alpha_1 - \beta_1$$

$$\text{suppose: } X_2 = \gamma_0 + \gamma_1 X_1 + W$$

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + V \\ &= (\beta_0 + \beta_2 \gamma_0) + (\beta_1 + \beta_2 \gamma_1) X_1 + (H + \beta_2 W) \end{aligned}$$

Omitted Variable Bias Formula

$$d_1 - \beta_1 = \beta_2 y_1$$

$\beta_2 > 0, y_1 > 0$, effect of x_1 will be overstated

Logarithms

- Consider a log-log specification

$$\log(Y) = \beta_0 + \beta_1 \log(X) + \epsilon$$

$$E[\Delta \log Y] = \beta_1 \Delta \log X$$

$$E\left[\frac{\Delta Y}{Y}\right] = \beta_1 \frac{\Delta X}{X}$$

β_1 is expected percent change in Y w/ 1% change in X

- consider Log-Linear Specification

$$\log Y = \beta_0 + \beta_1 X_1 + \epsilon$$

↓

$$E[100 \frac{\Delta Y}{Y}] = 100 \beta_1 \Delta X$$

$100 \beta_1$ is expected % change in Y w/ 1 unit increase X

Polynomials

Quadratic Specification

$$E[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2$$

$$\text{to } E[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2$$

Cubic Specification

$$E[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^3$$

$$\text{to } E[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^3$$

Dummy Variables

- we code educational attainment with Dummy Variables (Binary Variables)

$$\text{hsdrop} = \begin{cases} 1 & \text{if no high school} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{hscool} = \begin{cases} 1 & \text{if HS is highest} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{college} = \begin{cases} 1 & \text{if highest is college} \\ 0 & \text{otherwise} \end{cases}$$

Problem: Note

$$\text{hsdrop} + \text{hschool} + \text{college} + \text{morecol} = 1$$

↳ will result in Multicollinearity unless we exclude one dummy variable from Regression

Dummy variables: Example

Regression:

$$E[\text{Earnings}|\text{schooling}] = \beta_0 + \beta_1 \text{hschool} + \beta_2 \text{college} + \beta_3 \text{morecol}$$

"High School Dropout" is omitted

category or Base case

* The Intercept = mean of Base Case

$$\beta_0 = E[\text{Earnings}| \text{hsdrop}=1]$$

coefficients on the dummy variables are difference in mean of that group & base case

$$\beta_2 = E[\text{Earnings}| \text{college}] - E[\text{Earnings}| \text{hsdrop}]$$

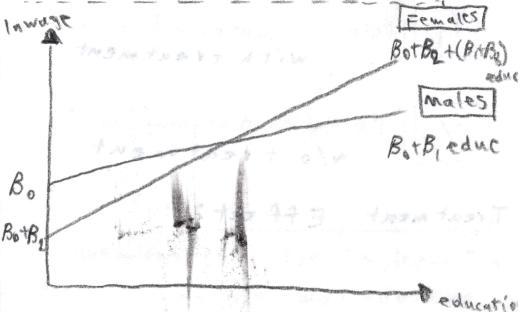
Interactions

- Let Female be a dummy variable for female worker

consider:

$$E[\text{In wage} | \text{female}, \text{education}] =$$

$$\beta_0 + \beta_1 (\text{education}) + \beta_2 (\text{female}) + \beta_3 (\text{education} \cdot \text{female})$$



- Here, Female * education is the interaction between females & education

Best Linear Predictors

$X^T B$ is best linear approximation to the regression function $E[Y|X]$

Non-Parametric Regression

use when:

- interested in shape
- We don't want to make functional form assumptions

Kernel Regression

- A Kernel Regression is a weighted average

$$\hat{E}[Y|X=x_0] = \sum_{i=1}^n w_i Y_i$$

where $w_i = \frac{K_h(x_i - x_0)}{\sum_j K_h(x_j - x_0)}$

$$w_i = \frac{K_h(x_i - x_0)}{\sum_j K_h(x_j - x_0)}$$

Function K_h is a Kernel

- $K_h(x_i - x_0)$ is Large if $|x_i - x_0|$ small
- $K_h(x_i - x_0)$ is small if $|x_i - x_0|$ large

Kernels

Usual Kernels:

1. Uniform



2. Gaussian



3. Epanechnikov

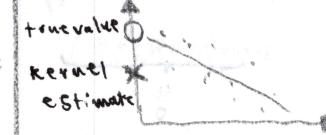


* h is the bandwidth: std. Dev of Kernel K_h

- If h is small, only observations close to x_0 get large weights
- If h is large, observations far from x_0 get larger weights

Kernel Regression: Boundary Bias

- consider x_0 at the boundary of support of x



* all data points far away, so its unreliable

Series Regression

- Fit a Polynomial of order K

$$\sum_{i=1}^n (Y_i - b_0 - b_1 X_i - b_2 X_i^2 - \dots - b_K X_i^K)^2$$

Local Linear Regression

- For each $X=x_0$, minimize

$$\sum_{i=1}^n K_h(x_i - x_0)(Y_i - b_0 - b_1 X_i)^2$$

Instrumental Variables

- say you collect log of daily prices log of quantity sold

We aim to estimate demand equation

$$\log q_t = \beta_0 + \beta_1 \log p_t + \epsilon_t$$

↳ β_1 = Price elasticity of Demand

Motivation

- Prices & Quantities determined by supply + demand

* Changes in H_t induce change in Price

$$\hookrightarrow \text{cov}(\log P_t, H_t) \neq 0 \quad \begin{cases} \text{OLS will} \\ \text{not work} \end{cases}$$

* We'll consider only supply shocks to estimate Demand function

Instrumental variable estimation

- We use bad weather as supply shock

$$\text{stormy}_t = \begin{cases} 1, \text{ if stormy day} \\ 0, \text{ otherwise} \end{cases}$$

$$\hookrightarrow \text{cov}(\text{stormy}_t, H_t) = 0 \quad \textcircled{1}$$

$$\hookrightarrow \text{cov}(\text{stormy}_t, P_t) \neq 0. \quad \textcircled{2}$$

- If $\textcircled{1}$ & $\textcircled{2}$ hold, stormy_t is an instrument for $\log P_t$

Instrumental Variable: Procedure

$$\text{cov}(\text{stormy}_t, \log P_t) = \beta_1, \text{cov}(\text{stormy}_t, \log P_t) + \text{cov}(\text{stormy}_t, H_t)$$

$$\hat{\beta}_1 = \frac{\text{cov}(\text{stormy}_t, \log P_t)}{\text{cov}(\text{stormy}_t, H_t)}$$

$$\hat{\beta}_0 = \log P_t - \hat{\beta}_1 \log H_t$$

Instrumental variables: General

$$Y_i = \beta_0 + \beta_1 X_i + H_i$$

- X_i is endogenous, when $\text{cov}(X_i, H_i) \neq 0$

- If there exists variable Z , where

$$\textcircled{1} \quad \text{cov}(Z, X) \neq 0$$

$$\textcircled{2} \quad \text{cov}(Z, H) = 0$$

* Z is a valid instrument for

endogenous variable X_i

$$\hat{\beta}_1 = \text{cov}(Y, Z) / \text{cov}(X, Z)$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

2 SLS Interpretation

- Instrumental Variable Estimator has 2-stage least squares interpretation

Purpose, Scope, & Examples

Program Evaluation assesses the causal effects of Public Policy interventions

Effect of:

1. Job training on earnings
2. Class size on test scores
3. Min. wage on employment

Causality w/ Potential Outcomes

Treatment

D_i : indicator of treatment intake of unit i

$$D_i = \begin{cases} 1, \text{ if unit } i \text{ receives treatment} \\ 0, \text{ otherwise} \end{cases}$$

Outcome

Y_i : observed outcome variable of interest for unit i

Potential Outcomes

Y_{1i} : Potential outcome for i with treatment

Y_{0i} : Potential outcome for i w/o treatment

Treatment Effect:

- Causal effect of treatment on outcome for unit i

$$Y_{1i} - Y_{0i}$$

Observed Outcomes

- Realized as

$$Y_i = Y_{1i} D_i + Y_{0i} (1 - D_i)$$

Fundamental Problem of Causal Inference

* cannot observe both potential outcomes

Identification Problem of Causal Inference

Problem: How to find $Y_{1i} - Y_{0i}$?

* homogeneity would solve,

- (Y_{0i}, Y_{1i}) constant over population or time

\hookrightarrow usually population is of heterogeneous nature

Stable Unit Treatment Value Assumption

Assumptions:

① treatment of i doesn't affect j

Quantities of Interest (Estimands)

ATE (average treatment effect)

$$d_{ATE} = E[Y_i - Y_0]$$

$ATET$ (average treatment effect on treated)

$$d_{ATET} = E[Y_i - Y_0 | D=1]$$

* Example given on Packet Selection Bias

$$E[Y|D=1] - E[Y|D=0]$$

$$= E[Y_i - Y_0 | D=1] + E[Y_0 | D=1] - E[Y_0 | D=0]$$

ATET

BIAS

Assignment Mechanism

- procedure that determines which units are selected for treatment intake

1. Random Assignment

2. Selection on Observables

3. Selection on unobservables

* Key Ideas *

- causality defined by Potential outcomes

* Observed Association is neither necessary nor sufficient for causation

- Estimation of causal effects starts w/ assignment mechanism