

PS699: Statistical Methods II

Quick Review of Fundamental Probability & Distributions, and Statistics, Estimation, & Inference

I. Basic Probability Review:

A. *Preamble:*

1. We can view the outcome of every aspect of the physical & social world as the product of some “random” “experiment”.
 - a) “Random” here means that prior to the outcome, some degree (from zero to infinite) uncertainty exists about what the outcome will be
 - (1) Now, zero & infinite uncertainty both, in some sense, strain the meaning of “uncertainty”. We are usually uninterested in either case:
 - (a) Zero, because there is nothing to explain; it just is. Infinite uncertainty is even little hard to grasp & probably not empirically relevant in vast majority of cases. We usually have some idea of range of possible outcomes if nothing else.
 - (b) Because of this, most texts exclude 0 and ∞ uncertainty from definition of “random”.
 - (2) Do remember, though, that something happens with probability 0 or with probability 1, either of which correspond to 0 uncertainty, is actually possible.
2. Another key term: ***Data-Generating Function (DGF)*** – the mechanism that produces these random outcomes.
 - a) Usually, in social science, DGF is outside analysts’ control.
 - b) One view of the object of our empirical-estimation exercises: to estimate the parameters of the (ideally, substantively & theoretically determined) DGF.

B. Probability Function:

So, let A represent the outcome of some random experiment, as in:
 $A \equiv$ some "event", a possible outcome of the experiment
then let $S \equiv$ set of all possible outcomes, called the Sample Space
and let $Pr(A) \equiv$ some function of A defined over the Sample Space

\Rightarrow Axiomatic Definition of a Probability Function (P.F.):

(Axiom)

Set Theory: ① $0 \leq Pr(A) \leq 1 \quad \forall A \in S$

• any event among permis.
has prob. b/w 0 & 1

② $Pr(S) = 1$

• some event among permis.
must happen.

③ $Pr(A \cup B) = Pr(A) + Pr(B) - Pr(A \cap B)$
 $\forall (A, B) \in S$

• this defines how we
manipulate probabilities
mathematically.

Notes:

on ① \rightarrow Notice that "nothing happens" is always an element of S b/c \emptyset (null set) always a subset of every other set, including any sample space.

on ③ \rightarrow if A & B disjoint, that is A & B have "nothing in common" or "no overlap"
then $Pr(A \cup B) = Pr(A) + Pr(B)$

C. *Random Variables*: assigning numeric values to outcomes probability functions gives definitions discrete and continuous probability functions

Now suppose for every possible outcome (ie, $\forall A \in \mathcal{S}$) we assign a numeric value, e.g. ① on a {roll of a die} there are {six possible faces that might come up}, we could assign to each {up face} the {numeric value on that face}

Annotations: random experiment (roll of a die), sample space (six possible faces), event or outcome (up face), random variable (numeric value)

② on a {coin toss} there are two possibilities, {heads or tails}, for each face that may come up, {H} or {T}, we assign a {value $H=1, T=0$ }

Annotations: random experiment (coin toss), sample space (heads or tails), event (H or T), random variable (value)

Random Variables: are mathematically defined thusly $X \equiv \begin{cases} 1 & \text{if coin heads up} \\ 0 & \text{if coin tails up} \end{cases}$

1. If outcomes finite or countably infinite, then R.V. is *discrete*.

define probability functions over (discrete) random variables axiomatically:

$f(x) \equiv \Pr(X=x)$

Annotations: probability function (p.f.) (f(x)), random variable (X), some placeholder for the specific values that X may take (x)

Axioms: ① $0 \leq \Pr(X=x) \leq 1$ or $0 \leq f(x) \leq 1$

② $\sum_{x_i \in \mathcal{S}} \Pr(X=x_i) = 1$ or $\sum_{x_i \in \mathcal{S}} f(x_i) = 1$

notice by these axioms that random variables are taken to be defined over disjoint events. That is, for any outcome only one value is assigned to X & for any value assigned to X only one outcome corresponds. (There is a "one-to-one" correspondence, or, identically, a function.)

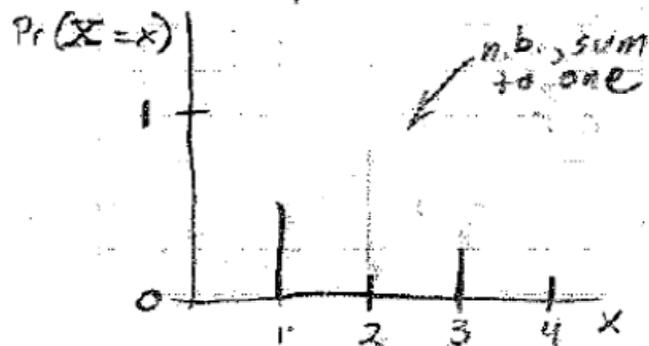
2. If outcomes neither *finite* nor *countably infinite*, then they are *infinitely divisible*, & the RV is *continuous*. \Rightarrow no way to sum probabilities of values \Rightarrow integrate, and the axiomatic definition of a **probability density function (pdf)** becomes:

$$\textcircled{1} \Pr(a \leq X \leq b) = \int_a^b f(x) dx \geq 0$$

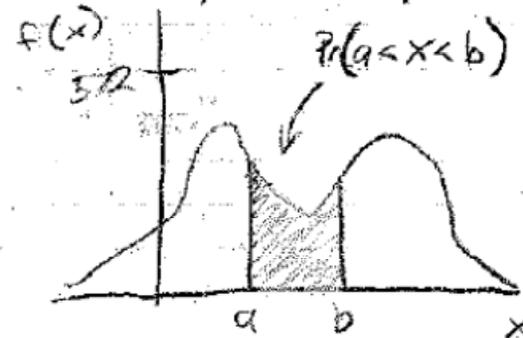
$$\textcircled{1a} \Pr(X = a) = \int_a^a f(x) dx = 0 \quad (\text{by definition of integrals})$$

$$\textcircled{2} \int_{-\infty}^{\infty} f(x) dx = 1$$

Probability Function



Probability Density Function



- n.b. $f(x)$ not $\Pr(X=x)$
- n.b. not nec. 0-1 on y axis
- n.b. can be any shape at all
- n.b. unbroken line (usu., not nec.)

- Notes:
- a) perfectly analogous to discrete case, except that probability at any specific point in the sample space is zero.
 - b) Thus $f(x)$ is not a probability function (pf) but ^{something else, we call it...} a probability density function (pdf)
 - c) Since $\Pr(X=a) = \emptyset$, $\Pr(a \leq X \leq b) = \Pr(a \leq X < b) = \Pr(a < X \leq b) = \Pr(a < X < b)$
 - d) Recall that integration is "scooping" up the area under a curve. So, to find the probability of some range of possible outcomes, say from a to b , we "scoop" the area under the pdf, $f(x)$, from a to b , i.e. we integrate from a to b .

D. Examples:

1. Discrete:

① Bernoulli (Binary Experiment)

- distributions have pf's or pdf's
- these have argument(s) they're called parameter(s) of the distribution.

what you need to know about the spec. example

- e.g.
- unemployed/not
 - coin toss
 - war/not war
 - individ. vote/not
 - gov't (un deficit/not
 - child is girl/boy
- } often, but not always, an $X \sim \text{not } X$ dichotomy

Bernoulli: $f(x) = \Pr(X=x)$ $x \in \{0, 1\}$

Bernoulli's parameter $\equiv \pi$

$$f(x) \equiv \Pr(X=x) = \begin{cases} \pi^x (1-\pi)^{1-x} & \text{for } x=0 \text{ or } x=1 \\ \emptyset & \text{for all else} \end{cases}$$

First Principles of Bernoulli (i.e., use it when what you're interested in follows or most closely follows these)

- ① only two possible outcomes
 - ② outcomes mutually exclusive (disjoint)
- ↔ but you can always redefine outcomes so that this is true

② Discrete Uniform (equiprobable) Distributions

$$f(x) \equiv \Pr(X=x) = \begin{cases} \frac{1}{N} & \text{for } x = 1, 2, \dots, N \\ \emptyset & \text{for all else} \end{cases}$$

parameter: N

eg. → coin toss
→ die roll
→ names drawn from hat
(which may approx some selection processes)

First Principles ① N disjoint events
② all equally likely

(Not used in our branches of stats much; used in formal theory a lot.)

③ Binomial Distribution:

$$f(x) \equiv \Pr(X=x) = \begin{cases} \frac{N!}{x!(N-x)!} \pi^x (1-\pi)^{N-x} & \text{for } x = 0, 1, 2, \dots, N \\ \emptyset & \text{for all else} \end{cases}$$

parameters:
 N, π

First Principles ① N Bernoulli trials
② Each trial has constant, same probability, π

(note: $N!/[x!(N-x)!]$ is “ N choose x ”, written $\binom{N}{x}$)

(4) **Poisson Distribution:** Limit as the number of trials $\rightarrow \infty$ at same time $\pi \rightarrow 0$ (e.g., because time-slice $\rightarrow 0$) such that $n\pi$, which is the mean of a binomial, is stable, we get continuous-time version of Binomial, the limiting distribution of the Binomial, known as the Poisson:

$$\Pr(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

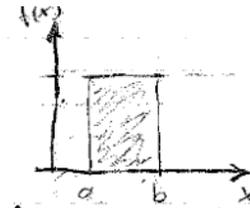
First-Principles: It's a count distribution. Each event is independent. Probability of event within observation period is λ .

2. Continuous:

① Continuous Uniform Distribution $[a, b]$

parameters:
 a, b -- the range
of possible
values

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for all } a < x < b \\ \phi & \text{else} \end{cases}$$



First Principles: ① fixed range of possible values (bounded)
 ② nothing "more likely" than anything else in this range

(even less/more commonly used in stats/theory than discrete)

② Normal (Gaussian) Function

$$f(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

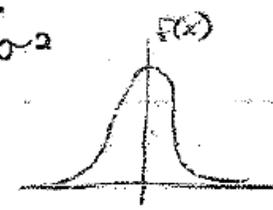
parameters:
 σ^2 → variance
 μ → mean

"the bell-shape"

• symmetric

• "thin-tailed"

no intuitive 1st principles, empirical (near) regularity



a) You'll see this *pdf* written other ways as well; for example:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left[\frac{(x-\mu)^2}{\sigma^2}\right]} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left[\frac{(x-\mu)}{\sigma}\right]^2} = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left\{\frac{-(x-\mu)^2}{2\sigma^2}\right\}$$

b) And, for reference, the **standard-normal** is the normal with mean 0 and s.d. 1, so:

$$x \sim N(0,1) \Leftrightarrow f(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\} = (2\pi)^{-0.5} e^{-\frac{x^2}{2}}$$

3. A squared *standard-normal* is *Chi-Squared*, χ_1^2 ; sum of n independent std-norms squared $\sim \chi_n^2$, sum of n independent $\chi_{n_i}^2 \sim \chi_{\sum_i n_i}^2$:

1. Chi-Squared Distribution

• if $Z \sim N(0, 1)$, then, letting $X \equiv Z^2$, $X \sim \chi_1^2$

(X is distributed Chi-Squared with 1 degree of freedom)

• $E(X)$ for any Chi-Squared r.v. X is its degrees of freedom

Aside: for any $X \sim N(\mu, \sigma^2)$, $Z = (X - \mu) / \sigma \sim N(0, 1)$

• if Z_1, Z_2, \dots, Z_n each dist. "independently" $N(0, 1)$, then $(\sum_{i=1}^n Z_i^2) \sim \chi_n^2$

• if X_1, X_2, \dots, X_n each dist. "independently" $\chi_{n_i}^2$, then $(\sum_{i=1}^n X_i) \sim \chi_{\sum n_i}^2$

Example: if the e_i from $y_i = X_i b + e_i$ distributed "independently" normal with variance σ_e^2 , then

$$\left(\sum_{i=1}^n \left(\frac{e_i}{\sigma_e} \right)^2 \right) \sim \chi_n^2$$

4. Ratio of two χ^2 , each divided by its degrees of freedom is F with those \circ free:

F Distribution: if X_1 & X_2 are "independent" chi-squared r.v.'s with n_1 & n_2 degrees of freedom respectively, then

$$F = \left(\frac{X_1/n_1}{X_2/n_2} \right) \sim F_{n_1, n_2} \rightarrow \text{"f is distributed F with } n_1 \text{ \& } n_2 \text{ degrees of freedom"}$$

Example: if e_i from $y_i = X_i b + e_i$ "independent" Normal $(0, \sigma_e^2)$ and e_i^* from $y_i^* = X_i^* b + e_i^*$ "independent" Normal $(0, \sigma_e^{*2})$ } "independent" from each other

$$\text{Then } \left\{ \frac{\left[\sum_{i=1}^{n_1} (e_i / \sigma_e)^2 \right] / n_1}{\left[\sum_{i=1}^{n_2} (e_i^* / \sigma_e^*)^2 \right] / n_2} \right\} \sim F_{n_1, n_2}$$

5. Standard normal, divided by square root of a χ^2 divided by its degrees of freedom is t with those deg's free:

t -Distribution (also called "student" t):

• if $Z \sim N(0, 1)$ & $X \sim \chi_n^2$ is "independent" of Z , then

$$\left(t \equiv \frac{Z}{\sqrt{X/n}} \right) \sim t_n \quad (\text{"t is distributed } t \text{ with } n \text{ deg. free"})$$

Example: if \hat{b} from $y_i = X_i b + e_i$ is $N(0, 1)$ and e_i is "independent" of \hat{b} , & also Normal $(0, 1)$

$$\text{then } \frac{\hat{b}}{\sqrt{Z e_i^2 / n}} \sim t_n$$

• it can be shown that $(t_n)^2 \sim F_{1, n}$ for any r.v. distributed t_n

6. Some useful relations & limiting distributions:

a) Square of a t_{n-k} is $F_{1,n-k}$: that is, $t_{n-k}^2 \sim F_{1,n-k}$

b) Two other useful limits: $\lim_{n \rightarrow \infty} t_{n-k}$ is $N(0,1)$; and $\lim_{n_2 \rightarrow \infty} (n_1 \times F_{n_1, n_2}) = \text{Chi-squared}[n_1]$

E. Cumulative Distribution Functions:

Cumulative Distribution Functions:
(cdf)

$$F(x) \equiv \sum_{X \leq x} f(x) \equiv \Pr(X \leq x) \quad (\text{discrete})$$

$$\Rightarrow f(x_i) = F(x_i) - F(x_{i-1})$$

continuous:

$$F(x) \equiv \Pr(X \leq x) \equiv \int_{-\infty}^x f(t) dt$$

$$\Rightarrow \frac{\partial F(x)}{\partial x} = f(x)$$

(a good test: you understand integration & differentiation well enough for our purposes if you can see why this is so)

math notes: since we want to integrate to x , we don't want to use x as the index (dx) by which we integrate, so we just use a different letter, say t

Further Implications:

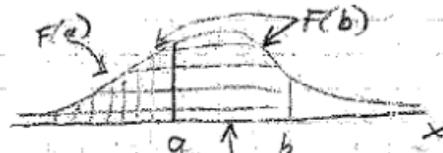
① $0 \leq F(x) \leq 1$ (by def. of prob. dens. function)

② if $x > y$, $F(x) \geq F(y)$ (by same def: x beyond y so "includes" it, $\Pr(X) \geq \emptyset$)

③ $F(+\infty) = 1$ (should be obvious) new stuff

④ $F(-\infty) = \emptyset$ (this too)

d ⑤ $\Pr(a \leq X \leq b) = F(b) - F(a)$



$$F(b) - F(a) = \int_a^b f(x) dx$$

F. Expectations (Means):

Measures of "central tendency". (trying to answer "where's this pf or pdf center?", "what's the 'middle' or 'center'?")

$$\begin{aligned} \text{"Mean of } X\text{"} &\equiv \mu_x \equiv E(X) \begin{cases} \equiv \sum_i X_i \cdot \text{pr}(X_i) \equiv \sum_i X_i \cdot f(X_i) & \text{discrete} \\ \equiv \int_{-\infty}^{\infty} X \cdot f(x) dx & \text{continuous} \end{cases} \\ &\equiv \text{"Expectation of } X\text{"} \end{aligned}$$

e.g. $E(X)$ where $X = \#$ dots on face of rolled die

$$E(X) = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 4 + \frac{1}{6} \cdot 5 + \frac{1}{6} \cdot 6 = \frac{21}{6} = 3\frac{3}{6} = 3.5$$

e.g. $E(X)$ where X continuous uniform on $0-1$:

$$E(X) = \int_{-\infty}^{\infty} X \cdot \frac{1}{(1-0)} dx = \int_0^1 X \cdot 1 dx = \int_0^1 X dx = \frac{1}{2} X^2 \Big|_0^1 = \frac{1}{2} - 0 = \frac{1}{2}$$

(actually ~~cont~~ uniform is always $\frac{1}{2}(b+a)$)

• Also median: "median of X " is X such that $\Pr(X \leq x) \geq \frac{1}{2}$ & $\Pr(X \geq x) \geq \frac{1}{2}$
↳ it's point where $\frac{1}{2}$ of pf or pdf on either side.
→ n.b. median not affected by extreme values of X , mean is

N.B. I say "mean of X " not "average of X ". In loose terms they are equivalent, but, more strictly, the mean is the true expectation of X in the pdf, the average is an estimate of the mean in some sample from that population. (More on this next week)

1. Median: not so commonly used in our branches of stats; again, more in theory
2. Mode: "most common value". That x at which $f(x)$ maximized, i.e., $\text{Argmax } f(x)$.

Expectation of Functions of X:

$$E(g(x)) = \sum_i g(x_i) \cdot f(x_i) \quad \swarrow \text{P.F.}$$

$$E(g(x)) = \int_{-\infty}^{\infty} g(x) f(x) dx \quad \begin{matrix} \uparrow \text{pdf} \\ \uparrow g(x) \end{matrix}$$

$$\Rightarrow \textcircled{1} E(a + bX) = a + bE(X) \quad \begin{matrix} \text{constants} \\ \downarrow \quad \downarrow \end{matrix}$$

$$\textcircled{2} E(a) = a$$

$$\textcircled{3} E(bX) = bE(X)$$

$$\textcircled{4} E(h(x) + g(x)) = E(h(x)) + E(g(x)) \quad \begin{matrix} \uparrow \text{any-old functions} \\ \uparrow \quad \uparrow \end{matrix}$$

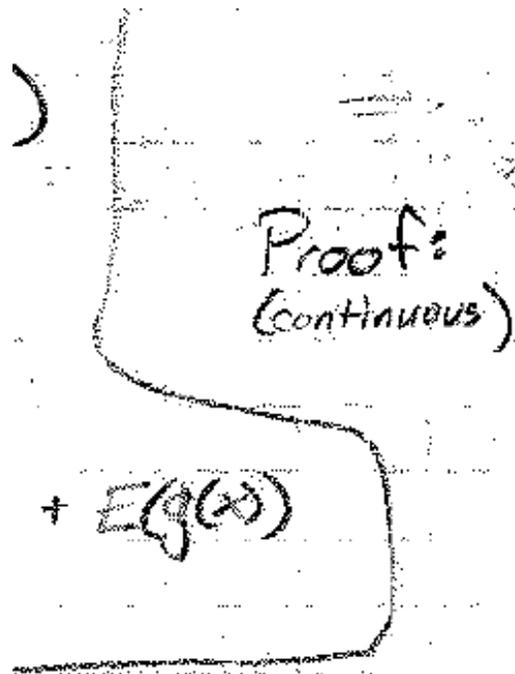
3. Rules for Working with Expectations:

- a) *Expectation of a sum is sum of the expectations.* (So, e.g., can move E in-&-out of Σ .)
- b) *Expectation operator slides through constants, snagging on random variables:*

$$E(a + bX + cZ + dXZ) = a + bE(X) + cE(Z) + dE(XZ)$$

- c) *Note: slides through constant coefficients, but cannot slide through product of 2 RV's*

4. How about a proof of this crucial property: $E(a+bX) = a + bE(X)$



Proof:
(continuous)

$$E(a+bX) = \int_{-\infty}^{\infty} (a+bX) f(x) dx$$

$$= \int_{-\infty}^{\infty} (af(x) + bx f(x)) dx$$

$$= \int_{-\infty}^{\infty} a f(x) dx + \int_{-\infty}^{\infty} bx f(x) dx$$

can do this by integral rules

$$= a \int_{-\infty}^{\infty} f(x) dx + b \int_{-\infty}^{\infty} x f(x) dx$$

can pull constants out by integral rules

= 1 by def of pdf

= E(X) by def of E(X)

$$= a \cdot 1 + bE(X)$$

$$= a + bE(X) \quad \square$$

5. Variance, crucial special case of $E(f(X))$:

VARIANCE: $V(X) \equiv \text{Var}(X) \equiv E[(X-\mu)^2] = \begin{cases} \sum (x_i - \mu)^2 \text{Pr}(X_i) & \text{discrete} \\ \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx & \text{continuous} \end{cases}$

$\equiv \sigma_x^2$

a) Some preliminary properties:

$$V(X) \equiv E[(X-\mu)^2] \equiv E(X^2 - 2X\mu + \mu^2) = E(X^2) - 2\mu E(X) + \mu^2 = E(X^2) - \mu^2$$

$$\textcircled{1} V(X) = E(X^2) - \mu^2$$

constants

$$= E(X^2) - [E(X)]^2$$

USE:
→ makes calc. variances

$$\textcircled{3} \text{Var}(g(x)) = \int_a^{\infty} (g(x) - E(g(x)))^2 f(x) dx$$

$$\approx [g'(\mu)]^2 \text{Var}(X)$$

$$\textcircled{2} \text{Var}(a + bX) = b^2 \text{Var}(X)$$

$$\textcircled{1} E(X^2) = \mu^2 + \sigma^2$$

$$\textcircled{2} \text{Var}(\text{any constant}) = 0$$

b) Note from (2): “constants neither vary nor *covary* [to be defined shortly...]”

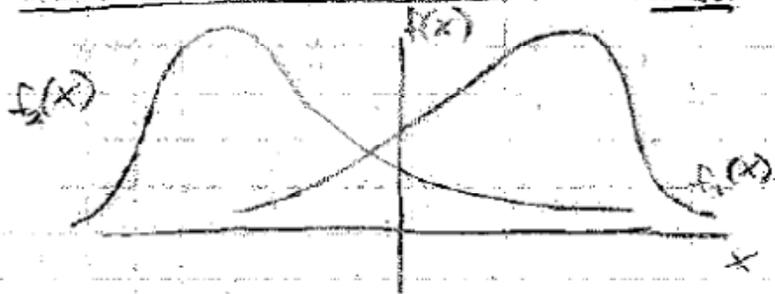
c) Note from (3): use Taylor-series linear-approximation for $V(\text{nonlinear } f(\text{RV}))$...

6. Mean Squared-Error, another crucial special case of $E(f(X))$:

$$MSE_{(x \text{ about } c)} \equiv E[(X-c)^2] = \underbrace{\sigma_X^2}_{\text{variance of } X} + \underbrace{(c-\mu)^2}_{\text{"bias" of } c \text{ relative to } \mu}$$

7. While at it, 3rd & 4th moments, Skew & Kurtosis also special cases $E(f(X))$:

Just some FYI definitions: skew: measure the symmetry of a distribution



$f_1(x)$ is skewed left (or negative)
 $f_2(x)$ is skewed right (or positive)
 → as Greene says, skew is direction of longer tail

$$\text{skew} \equiv E[(X-\mu)^3]$$

kurtosis: measures "how fat the tails" are (n.b. this not necessarily a simple relationship to variance)

e.g. t is more kurtotic than std normal.

$$\text{kurtosis} = E[(X-\mu)^4]$$

G. Joint (Multivariate) Distributions:

A. Suppose you have two R.V.'s, X_1 & X_2 (e.g. $X_1 = \begin{pmatrix} \text{on} \\ \text{male} \end{pmatrix}$ & $X_2 = \begin{pmatrix} \text{taller 6' or} \\ \text{shorter 6'} \end{pmatrix}$)

Defines: Joint (Bivariate) Distribution of X_1, X_2

or $X_1 = \begin{pmatrix} \text{GDP/capita} \end{pmatrix}$ & $X_2 = \begin{pmatrix} \text{Degree of Democ.} \end{pmatrix}$

$$\text{Prob}(a \leq X_1 \leq b, c \leq X_2 \leq d) = \begin{cases} \sum_{a \leq x_1 \leq b} \sum_{c \leq x_2 \leq d} f(x_1, x_2) & \text{discrete} \\ \int_a^b \int_c^d f(x_1, x_2) dx_2 dx_1 & \text{continuous} \end{cases}$$

Aside on doing "double integrals" & "double sums": just do them one at a time. First, sum or integrate over one variable, treating the other as a constant. Then, sum or integrate the other. (This generalizes to "multiple" sums & integrals) Changing the order of integration or summation makes no difference.

Axiomatic Definition:

- ① $f(x_1, x_2) \geq 0$
- ② $\begin{cases} \sum_{x_1, x_2} f(x_1, x_2) = 1 & \text{discrete} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) = 1 & \text{continuous} \end{cases}$

Cumulative Distribution Function (joint c.d.f.)

$$F(x_1, x_2) \equiv \text{Prob}(X_1 \leq a, X_2 \leq b) \quad \text{"means 'and'"} \quad \text{means "and"}$$

$$= \begin{cases} \sum_{x_1 \leq a} \sum_{x_2 \leq b} f(x_1, x_2) \\ \int_{-\infty}^a \int_{-\infty}^b f(x_1, x_2) dx_2 dx_1 \end{cases}$$

1. Marginal Distributions:

Marginal Distributions: (i.e., getting the univariate distributions back out of the bivariate distribution)

$$f_{X_1}(x_1) = \begin{cases} \sum_{x_2} f(x_1, x_2) \\ \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 \end{cases} \quad \text{d.v.v.} \quad f_{X_2}(x_2) = \begin{cases} \sum_{x_1} f(x_1, x_2) \\ \int_{-\infty}^{\infty} f(x_1, x_2) dx_1 \end{cases}$$

Intuition: A Snow-Covered Driveway

consider a driveway covered with snow. dimension
→ think of the length of the driveway as the X_1 dimension (length dimension, x-axis)
& the width of the driveway as the X_2 dimension (width, y-axis)

Then the height of the snow at any point on the driveway is as reflecting the "probability density" that any randomly given snowflake will have fallen at that (X_1, X_2) coordinate. In general, the height could vary from spot to spot of course.

First of all, we're defining snowflakes that fall on the driveway as the only relevant ones, so, by definition, the "probability density" at any point on the driveway is zero (no snow) or positive (some snow).

$$\Rightarrow f(x_1, x_2) \geq 0$$

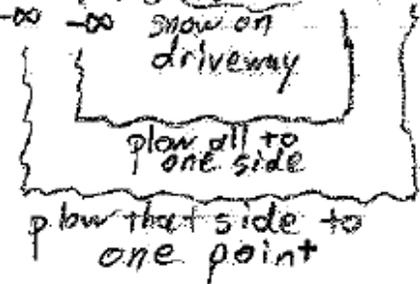
Since we've defined the driveway as the whole universe of relevant snow (seems that way shoveling it), any relevant snowflake must have fallen on the driveway somewhere.

$$\Rightarrow \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 dx_1 = 1$$

↳ or just integrate (shovel, snow-plow) over length & width of driveway.

So, if we took a shovel or a plow & plowed all the snow down to a line at the end or side of the driveway (that's the first integral), then turned the plow & plowed all the snow on that line into a point in the corner (that's the second integral after having done the first), we would have all the snow on that point. In probability, we call all : one.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 dx_1 = 1$$



⇒ all the probability, or 1.

A discrete example:

in % of population
= probability
randomly chosen
person will be in
that category

	h	
	6' and over	Under 6'
Male	18%	30%
Female	10%	42%
sum of rows	28%	72%

sum of (height) columns

48%
52%

$$\sum_{sex} \sum_{height} f(s, h) = \sum_{sex} (\text{sum of height})$$

$$= 48\% + 52\% = 100\%$$

$$= 1$$

no order or summation interdependence

Marginal Distributions: (i.e., getting the univariate distributions back out of the bivariate distribution)

$$f_{X_1}(x_1) = \begin{cases} \sum_{x_2} f(x_1, x_2) \\ \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 \end{cases}$$

$$\& \text{v.v. } f(x_2) = \begin{cases} \sum_{x_1} f(x_1, x_2) \\ \int_{-\infty}^{\infty} f(x_1, x_2) dx_1 \end{cases}$$

• back to the driveway analogy for intuition:

in general, snow falls all over the driveway, (x_1, x_2) coordinates all over. Suppose you wanted to know the probability (or proportionate amount) of snow falling at x_2 , not caring where it fell on relative to the x_1 axis. Say, proportionate amount in line with some tree next to the driveway. Well, get out your snow plow and plow the whole driveway (parallel to the street) so that it heaps up in a line along the side of the driveway (the line is perpendicular to the line from the tree of interest across the driveway). Wherever the snow was generally higher across the driveway will now be the highest point along the line you just plowed, won't it? (rhetorical) That height of that line of snow now represents the marginal (i.e. univariate) distribution of x_1 as it gives the "probability density" that x_1 coordinate was that point whatever x_2 was. (i.e., the "accumulated" height of snow along the line at x_1)

Example: in the sex & height table above, height = $\begin{cases} \text{over six foot w/ } p = 28\% \\ \text{under six foot w/ } p = 72\% \end{cases}$ is the marginal distribution of height. It's obtained by summing rows i.e. summing over sexes. — marginal.

2. Expectations of RV's & Functions of RV's in Multivariate Distributions:

Expectations From Joint p.d.f.'s & p.f.'s: from the definition just given, we can show a whole bunch of other useful things. ⑤

$$1) \quad E[X_1] = \sum_{x_1} x_1 f_{x_1}(x_1) \\ = \sum_{x_1} x_1 \left(\sum_{x_2} f(x_1, x_2) \right)$$

$$= \sum_{x_1} \sum_{x_2} x_1 f(x_1, x_2)$$

• from the definition of the marginal dist., f_{x_1}

• x_1 is not changed as we sum over x_2 , i.e., relative to summing over x_2 it is a constant, thus we can move it in or out of summing over x_2 at will.

• The Continuous Analogue is

$$E(X_1) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1, x_2) dx_2 dx_1$$

$$2) \quad V(X_1) \equiv E[(X_1 - \mu_{x_1})^2] \\ = \sum_{x_1} (x_1 - \mu_{x_1})^2 f_{x_1}(x_1)$$

$$= \sum_{x_1} \sum_{x_2} (x_1 - \mu_{x_1})^2 f(x_1, x_2) \\ = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_1 - \mu_{x_1})^2 f(x_1, x_2) dx_2 dx_1$$

• logic for this exactly as above

• by analogy

notational note: I may from time to time write it this way if both (or all) integrals over same range. If I ever leave it off, or write $\int_{x_1} I$ mean $-\infty$ to ∞ (not indefinite integral) or over the range of x_1 .

3. Covariance & Correlation (special cases of $E(f(X))$ in multivariate dists):

Covariance & Correlation: Preliminary Point: $E[g(x_1, x_2)] = \sum_{x_1} \sum_{x_2} g(x_1, x_2) f(x_1, x_2)$ (any function, including $g(x_1, x_2) = (x_1, x_2)$ joint p.f.)

Covariance $(X_1, X_2) \equiv \text{Cov}(X_1, X_2) \equiv C(X_1, X_2) \equiv E[(X_1 - \mu_{X_1})(X_2 - \mu_{X_2})]$

$= E[X_1 X_2 - \mu_{X_1} X_2 - X_1 \mu_{X_2} + \mu_{X_1} \mu_{X_2}]$

$= E(X_1 X_2) - \mu_{X_1} E(X_2) - E(X_1) \mu_{X_2} + \mu_{X_1} \mu_{X_2}$

$= E(X_1 X_2) - \mu_{X_1} \mu_{X_2} - \mu_{X_1} \mu_{X_2} + \mu_{X_1} \mu_{X_2}$

$C(X_1, X_2) = E(X_1 X_2) - \mu_{X_1} \mu_{X_2} \equiv \sigma_{X_1 X_2}$

Correlation $(X_1, X_2) \equiv \text{Corr}(X_1, X_2) \equiv \rho_{X_1, X_2}$

$\equiv \frac{\sigma_{X_1 X_2}}{(\sigma_{X_1} \sigma_{X_2})}$ (where σ_{X_i} is std. dev. X_i , etc.)

intuition: when X_1 is high, does X_2 tend to be high? if so, then this tends to be large. If when X_1 is high, does X_2 tend to be low? Then this product will be large & negative. If there's no relation, the product will tend to be zero. Question: if X_1 & X_2 vectors of observ. on R.V.'s & X_1 & X_2 are orthogonal, $(X_1, X_2) = 0$, what's the covariance of the R.V.'s X_1 & X_2 ?

- a) Both covariance and correlation are measures of simple (i.e., linear) association, with larger positive (negative) values indicating greater positive (negative) association, smaller indicating smaller, and 0 indicating no (simple, linear) association.
- b) Covariance, however, is sensitive to the scale of the variables; correlation nets those scales to produce a single measure $-1 \leq \rho \leq 1$, with -1 (+1) indicating perfect negative (positive) association and 0 indicating no association.

4. Conditional Distributions & Expectations (very special cases of $E(f(X))$...):

Conditional Distributions

$$f(x_2 | x_1) \equiv f(x_1, x_2) / f_{x_1}(x_1)$$

conditional distribution
"f of x_2 , conditional on x_1 ."
"f of x_2 given x_1 ."

writing this $f(x_1, x_2) = f(x_2 | x_1) f_{x_1}(x_1)$
also has important uses.

(definition works for discrete or continuous; in fact also for set theory analog)

$$P(A|B) = P(A \cap B) / P(B)$$

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

(This also goes by "Bayes' Law")

$\forall f_{x_1}(x_1) \neq 0$
obvious since you can't be "given" impossible information (given x_1 when $p(x_1) = 0$)

So, within any joint distribution, $f(x_1, x_2)$ there are 2 marginal distributions, $f_{x_1}(x_1)$ & $f_{x_2}(x_2)$, (or n of them in multivariate), & two conditional distributions, $f(x_1 | x_2)$ & $f(x_2 | x_1)$, (or a whole boatload (I think $N(N!)$) in the multivariate case).

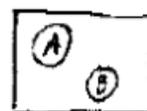
Conditional Expectations

We may very well be interested in the expectation of x_2 given x_1 , or renaming the variables for obvious reasons, y given x .

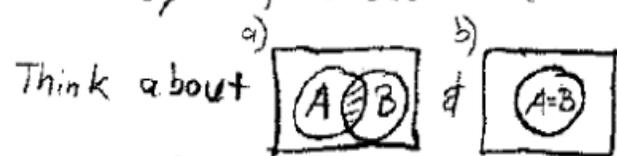
$$E(y|x) = \begin{cases} \sum_y y \cdot f(y|x) \\ \int_{-\infty}^{\infty} y \cdot f(y|x) dy \end{cases}$$

$(E(y|x)$ sometimes written $\mu_{y|x}$)
Greene also calls this "regression of y on x "

analogous



$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0}{P(B)} = 0$
"p of A given B" is obviously zero, so check!



in b) $P(A \text{ given } B) = 1$
check = $P(A \cap B) / P(B) = \frac{P(B)}{P(B)}$

in a) $\frac{P(A \cap B)}{P(B)}$ is the shaded as a fraction of B which reasons out well enough: once we're given that we're in B, we can think of B as the whole sample space \rightarrow

Example: Go back to our sex & height table.

that's right, but note not necessarily "linear regression of y on x"

in the whole sample space $P(A|B)$ is just the part intersects as fractions

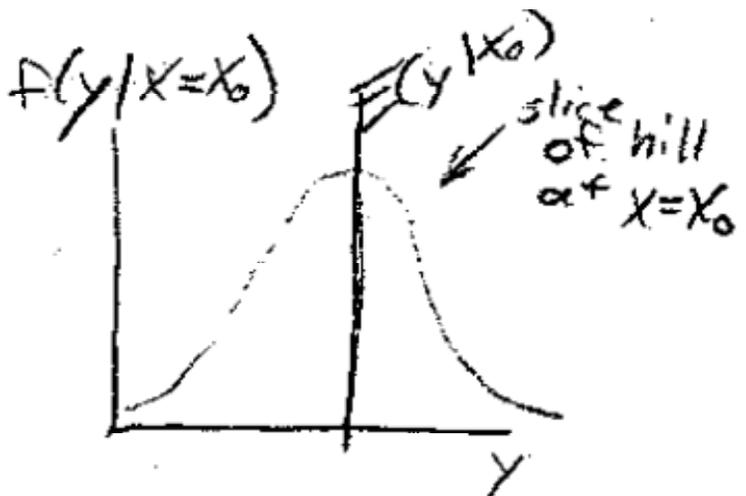
- The conditional distribution of (sex | height $\geq 6'$) is $f(y|x) = \begin{cases} \Pr(\text{male} \& 6') / \Pr(6') = .18 / .28 = 9/14 \\ \Pr(\text{female} \& 6') / \Pr(6') = .10 / .28 = 5/14 \end{cases}$
- The conditional expectation, calling female = 1 & male = 0,

is
$$E(y|x) = \sum_y y \cdot f(y|x) = 0 \cdot \frac{9}{14} + 1 \cdot \frac{5}{14} = \frac{5}{14}$$

$f(y|x=x_0)$ $E(y|x_0)$ slice

Harder to draw for continuous, but suppose...

$14 + 1 \cdot \frac{9}{14} = \frac{5}{14}$



These are "hills" are coming out of the page, darker is higher, higher represents "more probable" i.e. larger probability density a slice of the hill at x_0 , say, would be 2 dimensional & could be drawn as in the second chart. That's $f(y|x)$ $E(y|x_0)$ is just the mean of that slice.

Perhaps the most important example is the *expected squared, conditional deviation* function or the *conditional-variance* function, i.e., the *skedastic* function:

$$V(y|x) \equiv E[(y - E(y|x))^2 | X] = \begin{cases} \sum_y [y - E(y|x)]^2 f(y|x) \\ \int_{-\infty}^{\infty} (y - E(y|x))^2 f(y|x) dy \end{cases}$$

writing it out, we can simplify:

$$\begin{aligned} V(y|x) &\equiv E[(y - E(y|x))^2 | X] = E[y^2 - 2yE(y|x) + \{E(y|x)\}^2 | X] \\ &= E[y^2 | X] - 2E[yE(y|x) | X] + E[\{E(y|x)\}^2 | X] \\ &= E[y^2 | X] - 2E[y \mu_{y|x} | X] + [E(y|x)]^2 \\ &= E[y^2 | X] - 2\mu_{y|x} E(y|x) + \mu_{y|x}^2 \end{aligned}$$

(Expect. of sum is sum of expects.)
n.b., "givens" get distributed

law of iterated expectations, "giving X" twice is redundant

$\Rightarrow V(y|x) = E(y^2|x) - [E(y|x)]^2$

Conditional Variance: $V(y|x) \equiv E[\{y - E(y|x)\}^2 | X]$

$$\begin{aligned} &\equiv \sum_y (y - \mu_{y|x})^2 \cdot f(y|x) \quad \text{or} \quad \int_{-\infty}^{\infty} (y - \mu_{y|x})^2 f(y|x) dy \\ &= E(y^2|x) - \mu_{y|x}^2 \end{aligned}$$

• Returning to the (beautiful) contour map of the "hills" & the slice from that map: the conditional variance of y given $X = X_0$ is the variance of $f(y|x)$ which is the variance of that slice of the "hill" i.e., a measure of how "spread" the distribution is, how "broad" or "fat" the "hill" is at that point. In general, the cond. var. depends on which x you give bc cross-sectional (conditional) shape of the hill is going to differ by where you slice it, unless the hill is very regularly shaped. No "true hill" would, for example, have identically shaped cross-sections.

FYI: The conditional-variance function is called the skedastic function. A (joint) distribution with identical $V(y|x) \forall x$, i.e. the same skedastic function for any x is called homoskedastic. If $V(y|x)$ varies by x we give, then the (joint) distribution is heteroskedastic.

- NB: Expectations of $y|x$: $E(y|x)$, $V(y|x)$ are generally going to be functions of x or simply constants. Thus, they "use" x to "describe" (give central tendency & spread of) y .
- FYI: As thinking about the hill-crosssection analogy may make clear, there are conditional skewness, conditional kurtoses, etc. (i.e., conditional moments).
- Alternative Analogy for Conditional Distributions: Think of a sheet hovering in a room over air jets blowing at various strengths.
 - The strength with which a particular part of the sheet is being "blown" is actually the force of frequency with which it is hit by particles in the air stream. Thus the height of the sheet is actually a pretty good analogy for a joint p.d.f. The coordinates, x & y , give the layout of the room, & the height of the sheet at (x, y) is $f(x, y)$, the p.d.f. (a strength with which it's being "blown").
 - The conditional distribution $f(y|x=x_0)$, say, would be by the thread running through the sheet at x_0 .
 - $f(x|y=y_0)$ is analogously, the thread running the other way at y_0 .

5. Summary: Key Properties of Expectations, (Co)Variances, & Conditioning

Ⓐ Expectations

1) Laws of Iterated Expectations

x, y, z
= R.V.'s

$$a) E(E(x)) = E(x)$$

a, b, c
= constants

$$b) E(E(y|x)) = E(y)$$

2) Expectation of a sum is sum of the expectations

$$E\left(\sum_{i=1}^n g(x)\right) = \sum_{i=1}^n E(g(x))$$

3) Expectation of a constant is the constant

$$a) E(a) = a \quad b) E(x|x) = x$$

to "give" something renders it constant
(if "conditions on that info")

4) Rob's Rule of Expectations:

"E slides on through constants, snagging around random variables"

$$a) E(a + bX + cY) = a + bE(X) + cE(Y)$$

$$b) E\left\{ (X'X)^{-1} X' [ee'] X (X'X)^{-1} \right\} = (X'X)^{-1} X' E(ee') X (X'X)^{-1}$$

↑ constant matrix ↑ random vector

(B) Conditioning:

1) Conditioning is distributive:

$$E(a + bX + cY | X) = E(a | X) + bE(X | X) + cE(Y | X) \\ = a + bX + cE(Y | X)$$

2) Conditioning More than Once is Redundant:

$$E[E(Y | X) | X] = E(Y | X) \quad \leftarrow \text{(on same info!)}$$

3) Conditioning on X renders X constant (see ABb)

4) Conditioning & Expectation sort of undo each other (see ABb, ABc)

Rob's Rules for Variances & Covariances

a) Constants neither vary nor covary.

$$V(a) = 0 \quad C(a, X) = 0$$

b) "Variance slides through sums, eliminating constants, snagging on random variables and squaring their coefficients, while spitting out 2 times all the covariances times their associated coefficients."

$$\rightarrow V(a + bX + cY) = b^2 V(X) + c^2 V(Y) + 2bc \text{Cov}(X, Y)$$

→ Notice that in Matrix notation, with X a constant vector and b a random vector, this becomes just:

$$V(\underbrace{X}_{\sim}' \underbrace{b}_{\sim}) = \underbrace{X}_{\sim}' \underbrace{V(b)}_{\sim} \underbrace{X}_{\sim}$$

V-cov mat

↑
"X squared" in matrix-land

→ Special Cases of the Rule:

- $V(X + Y) = V(X) + V(Y) + 2 \times \text{Cov}(X, Y)$
- $V(X - Y) = V(X) + V(Y) - 2 \times \text{Cov}(X, Y)$

↑ $(-1)^2$ ↑ $2 \times (-1)$

c) The Analogous for Covariances (becomes a kind of FOIL):

$$\text{Cov}(a + bX + cY, d + eW + fV) =$$

$$be \text{Cov}(X, W) + bf \text{Cov}(X, V) + ce \text{Cov}(Y, W) + cf \text{Cov}(Y, V)$$

⑤ Properties of Note Combining the Above:

1) Any R.V. y can be written: $y = E(y|x) + (y - E(y|x))$
 $\equiv E(y|x) + \varepsilon$

2) Properties that follow:

a) $Cov(X, Y) = Cov(X, E(Y|X))$

$\equiv \hat{y} + \varepsilon$
 n.b., not necessarily linear \hat{y} , & true $E(y|x)$ not literally \hat{y}

b) $V(Y) = E(V(Y|X)) + V(E(Y|X))$

$= \hat{\sigma}_\varepsilon^2 + V(\hat{y})$ } "variance decomposition"
 $\Rightarrow TSS = RSS + ESS = ESS + RSS$
 $= SSR + SSE = SSE + SSR$
 (Labels: residual, explained, error, regression)

Note: $\Rightarrow V(\hat{y}) \leq V(y)$ (and also that $V(\varepsilon) \leq V(y)$ though that less notable)

I.D.2. (cont.) Properties that follow from $E(y) = E(y|x) + (y - E(y|x))$
 $= \hat{y} + \varepsilon$

c) $E(XY) = E(X \cdot E(Y|X))$

d) $E(\varepsilon) = E(\varepsilon|x) = E(\varepsilon|y) = 0$

- i) $E(\varepsilon) = E(y - E(y|x)) = E(y) - E(E(y|x)) = E(y) - E(y) = 0$
- ii) $E(\varepsilon|x) = E(y - E(y|x)|x) = E(y|x) - E(E(y|x)|x) = E(y|x) - E(y|x) = 0$
- iii) $E(\varepsilon|y) = E(y - E(y|x)|y) = E(y|y) - E(E(y|x)|y) = y - y = 0$

e) $Cov(X, \varepsilon) = 0$; in fact, $Cov(h(x), \varepsilon) = 0 \forall h(\cdot)$

$\hookrightarrow E(X\varepsilon) - E(X)E(\varepsilon) = E(Xy - XE(y|x)) - E(X) \cdot 0$
 $= E(Xy) - E(XE(y|x))$
 $= E(Xy) - E(Xy) = 0$

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

$$\hookrightarrow V(Y - E(Y|X)) = V(Y) + V(E(Y|X)) - 2\text{Cov}(Y, E(Y|X))$$

$$= V(Y) + V(E(Y|X)) - 2V(E(Y|X))$$

$$= V(Y) - V(E(Y|X)) \quad \square$$

I. D.3. For Special Case where $E(Y|X)$ is linear, ^{using variance-decomposition result} i.e. $E(Y|X) = a + bX$

$$a) a = E(Y) - bE(X)$$

$$b) b = \text{Cov}(X, Y) / V(X) \quad \text{proof: } \begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, E(Y|X)) \\ &= \text{Cov}(X, a + bX) \\ &= bV(X) \\ \Rightarrow b &= \text{Cov}(X, Y) / V(X) \quad \square \end{aligned}$$

$$c) V(Y) = E(V(Y|X)) + V(E(Y|X))$$

$$= E(V(Y - a - bX)) + V(a + bX)$$

$$= E(V(E)) + b^2 V(X) \quad \begin{array}{l} \text{recall that "giving" } X, \text{ fixes it} \\ \text{if constants don't} \\ \text{vary or covary} \end{array}$$

$$d) \rho_{X,Y} = \sqrt{b_{xy} b_{yx}} = \sqrt{\frac{\text{Cov}(X, Y)}{V(Y)} \times \frac{\text{Cov}(Y, X)}{V(X)}} = \sqrt{\frac{(\text{Cov}(X, Y))^2}{V(Y)V(X)}}$$

$$e) E(Y) = a + bE(X); \text{ i.e., } \bar{y} = a + b\bar{x} \quad = \frac{\text{Cov}(X, Y)}{\text{sd}(X)\text{sd}(Y)} \quad \square$$

$$f) \text{ "Coefficient of Determination," } R^2; = \{\rho(Y, \hat{Y})\}^2$$

H. Independence:

① Stochastic (Statistical) Independence $f(y|x) = f(y)$

$$f(x,y) = f(x)f(y)$$

$$f(x|y) = f(x)$$

proofs:

$$f(y|x) \equiv \frac{f(x,y)}{f(x)}$$

so, if $f(y|x) = f(y)$, then $f(x,y) = f(x)f(y)$

uni-directional

$$f(x|y) \equiv \frac{f(x,y)}{f(y)}$$

so, if $f(x|y) = f(x)$, then $f(x,y) = f(x)f(y)$

if $f(x,y) = f(x)f(y)$, then $f(x) = \frac{f(x,y)}{f(y)} \equiv f(x|y)$

& $f(y) = \frac{f(x,y)}{f(x)} \equiv f(y|x)$

knowing X doesn't tell you one darn thing about the distribution of Y & v.v. Not mean, not variance, nothing

② Mean Independence: $E(Y|X) = E(Y)$ corollary if $E(Y|X) = E(Y)$

$$E(Y|X) = E(Y) \not\Rightarrow E(X|Y) = E(X)$$

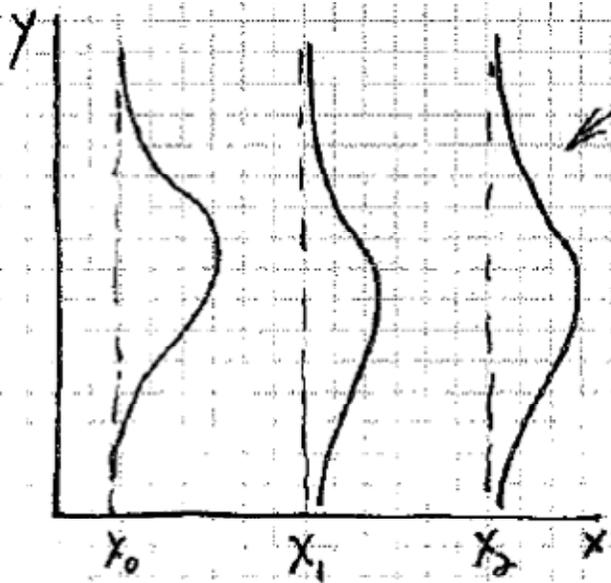
the $E(Y|h(x)) = E(Y) \forall h(\cdot)$

knowing X tells you $E(Y)$, the mean of Y ; may not tell you \sqrt{Y} or other stuff about the distribution of Y

③ "Linear Independence" or Uncorrelatedness
(clarify relation to 'linear independent columns')

$$\text{Cov}(X, Y) = 0 \iff \text{linear independence}$$

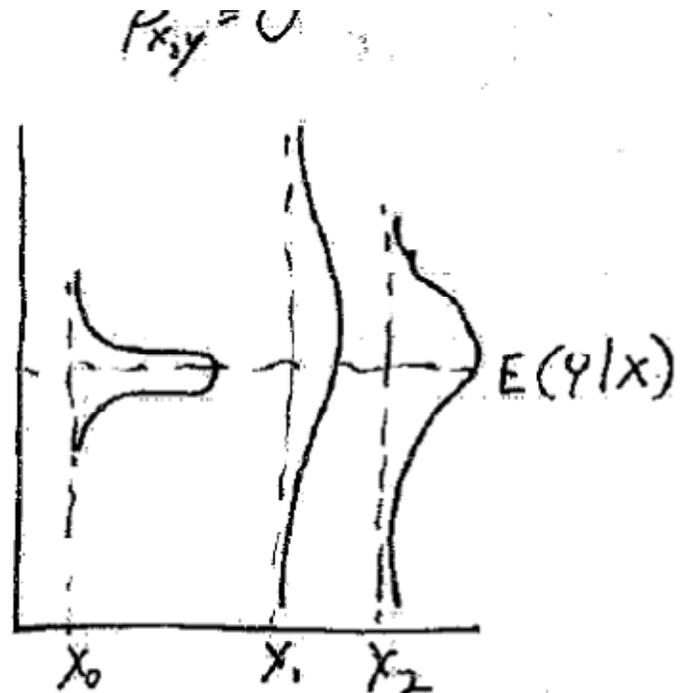
$$\rho_{X,Y} = 0$$



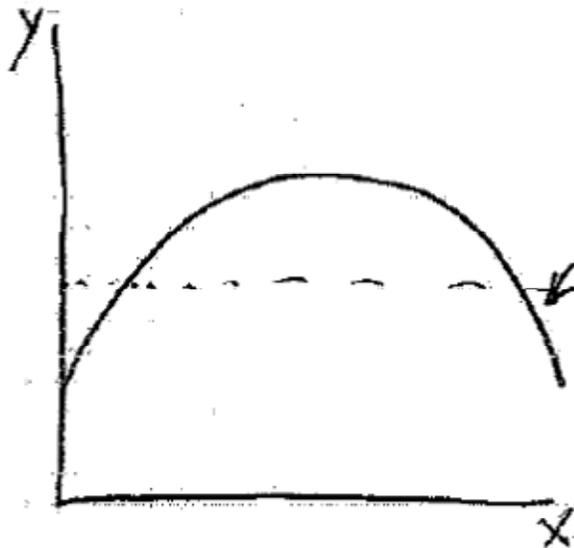
Stochastic

same ~~shape~~ hills exactly proportional

Similarly, proportional from other direction also



mean but not stoch



$E(y|x)$
 $E_{stoch}(y|x)$

} linear but not mean

$$E(\underline{X}) = \begin{bmatrix} E(X_{i_1}) & \dots & E(X_{i_n}) \\ \vdots & & \vdots \\ E(X_{j_1}) & \dots & E(X_{j_n}) \end{bmatrix}$$

$$E(\underline{\varepsilon}) = \begin{bmatrix} E(\varepsilon_1) \\ \vdots \\ E(\varepsilon_n) \end{bmatrix}$$

$$V(\underline{\varepsilon}) = E[(\underline{\varepsilon} - \underline{\mu})(\underline{\varepsilon} - \underline{\mu})'] = \begin{bmatrix} \sigma_1^2 & \dots & \sigma_n^2 \\ \vdots & & \vdots \\ \sigma_1^2 & \dots & \sigma_n^2 \end{bmatrix}$$

or Ω

Symmetric

$$V(\underline{b}) = \begin{bmatrix} V(b_1) & \dots & C(b_1, b_k) \\ \vdots & & \vdots \\ C(b_k, b_1) & \dots & V(b_k) \end{bmatrix}$$

confusion in notes regarding $V(\underline{X})$

Special Case: if X_1, X_2, \dots, X_n multivariate normal,
then linear \Rightarrow mean \Rightarrow stochastic

Longer-hand demonstration of $V(\text{vector}) = \text{matrix}$:

$$V(\underline{\epsilon}) = E[(\underline{\epsilon} - \underline{\mu})(\underline{\epsilon} - \underline{\mu})'] = E \begin{bmatrix} (\epsilon_1 - \mu_1)^2 & \dots & (\epsilon_1 - \mu_1)(\epsilon_2 - \mu_2) \\ \vdots & (\epsilon_2 - \mu_2)^2 & \vdots \\ (\epsilon_1 - \mu_1)(\epsilon_2 - \mu_2) & \dots & (\epsilon_n - \mu_n)^2 \end{bmatrix}$$

$$= E \begin{bmatrix} \epsilon_1 \epsilon_1 & \epsilon_1 \epsilon_2 & \dots & \epsilon_1 \epsilon_n \\ \epsilon_2 \epsilon_1 & \epsilon_2 \epsilon_2 & & \\ \epsilon_3 \epsilon_1 & & & \\ \vdots & & & \\ \epsilon_n \epsilon_1 & \dots & & \end{bmatrix}$$

$$= E \begin{bmatrix} \epsilon_1^2 & \dots & \epsilon_1 \epsilon_n \\ \vdots & \epsilon_2^2 & \vdots \\ \vdots & \vdots & \ddots \\ \epsilon_n \epsilon_1 & \dots & \epsilon_n^2 \end{bmatrix}$$

$$\Rightarrow \underline{\Sigma} = \underline{\Omega} = \begin{bmatrix} E(\epsilon_1^2) & \dots & E(\epsilon_1 \epsilon_n) \\ \vdots & \ddots & \vdots \\ E(\epsilon_n \epsilon_1) & \dots & E(\epsilon_n^2) \end{bmatrix} = \begin{bmatrix} V(\epsilon_1) & \dots & C(\epsilon_1, \epsilon_n) \\ \vdots & \ddots & \vdots \\ C(\epsilon_n, \epsilon_1) & \dots & V(\epsilon_n) \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \dots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \dots & \sigma_n^2 \end{bmatrix}$$

CRUCIAL IMPLICATION of INDEPENDENCE:

Given independence, joint distribution is just product of the marginal distributions:

Since in SRS, each $X_i \sim$ independently $f(x_i, \theta)$, the joint distribution of the X_i 's is

$$f(\underset{\sim}{X}, \underset{\sim}{\theta}) = f(x_1, \theta) f(x_2, \theta) f(x_3, \theta) \dots$$

vector of \rightarrow
the X_i 's or
matrix of the
vectors $\underset{\sim}{X}_i$

$$= \prod_{i=1}^n f(x_i, \theta)$$

II. Estimation & Statistical Inference

A. Preamble:

We will think of outcomes or events in the world as having been generated by some underlying (and usually unobservable) probabilistic process.

A. Thus, we will think of outcomes as having been generated by some probability function or probability density function.

1. E.g., we might think of election outcomes, ^(in 2-way competitions) as having been produced by some Bernoulli probability function.
2. E.g., we might similarly think of how much the two candidates spend on the campaign as coming from p.d.f.'s

- a) The sum of their spending might be considered $\text{Normal}(\mu, \sigma^2)$
- b) Each candidate's spending might be considered to have been drawn from a bivariate normal $N(\mu, \Sigma)$

B. This does not mean that we think of the world as random in the sense of "completely unpredictable" or "having nothing systematic about it"

1. For one thing, the statement that the outcome in question comes from some particular type of p.f. or p.d.f. is itself a statement of "probabilistic regularity"
2. For another thing, we usually go on to say that we know or think we know something about what determines the parameters of the distribution
 - a) We may think, for example, that π_i , probability the incumbent wins is some function of the amount of pork s/he brought to her/his district
 - b) Total spending we might think to be a function of how close to 50-50 the partisanship of voters in the district, for example

C. We might, then, wish to use the outcomes that actually occurred in some real-world sample of like events to infer something about the (parameters of the) underlying distribution.

D. This is the whole point: we are interested in inferring from what we observe how the world actually works; the underlying systematic elements of the political $\&$, more broadly, the socio-political-economic world.

1. This course is about how best to use the information available from what has happened.
2. "Best" here refers to the "science" of accurate inference; it is decidedly an applied science as we are interested in it

III. Samples, Sampling Theory, & Sample Statistics

A. (Simple) Random Sample (SRS): n observations, $X_1, X_2, X_3, \dots, X_n$ drawn independently from the same population (i.e. from the universe of possibilities as described by a p.f. or p.d.f.) is a S.R.S. from that p.f. or p.d.f.:

$\underbrace{X_1, \dots, X_n}_{\sim}$ a S.R.S. from $f(\underbrace{X_i}_{\sim}, \underbrace{\theta}_{\sim})$

↑
may be an observation on a single R.V. or on a random vector or even a random matrix

↑
the p.f. or p.d.f.

↑
the parameter, or parameters of that p.d.f.

e.g. X_1, \dots, X_n a S.R.S. from $N(\mu, \sigma^2)$

\equiv each X_i a draw from an independent $N(\mu, \sigma^2)$

$\rightarrow \theta$

1) In practice, our samples will be in the form of a

a) cross-section: a set of observations on various units
(CS) (individuals, states, countries, governments, etc.)
at some time, t .

b) time-series: a set of observations on a single unit
(TS) across some number of (usually evenly spaced)
time periods, $t=1, 2, 3, \dots, T$

c) time-series-cross-section: the combination of (a) & (b) -- a set
(TSCS) of cross-sectional observations across some
number of time periods, or equivalently a set
of time-series observations across some
cross-section of units.

FYI: (i) a TSCS of data with (many) more CS than TS is
usually called panel data.

(ii) a TSCS of data with (many) more TS than CS is
usually called pooled data.

→ this nomenclature is not all that terribly adhered to though

2) Since in SRS, each $X_i \sim$ independently $f(X_i, \theta)$, the joint
distribution of the X_i 's is

$$f(\underset{\sim}{X}, \underset{\sim}{\theta}) = f(x_1, \underset{\sim}{\theta}) f(x_2, \underset{\sim}{\theta}) f(x_3, \underset{\sim}{\theta}) \dots$$

vector of \rightarrow
the X_i 's or
matrix of the
vectors $\underset{\sim}{X}_i$

$$= \prod_{i=1}^n f(x_i, \underset{\sim}{\theta})$$

B) (Sample) Descriptive Statistics

1) A statistic is any function computed from the data.

• of course, we usually restrict our attention to functions of the data that are informative about the distribution of the data (and hopefully \therefore , by some inference process, informative about the underlying population).

2) Some common "descriptive" statistics (i.e. they purport to help describe or summarize the entire sample distribution & , again by inference, the population distribution)

III. B. 2. Descriptive Statistics



a) First, generically, a statistic, δ is some function of the data
$$\delta = g(x_1, x_2, \dots, x_n)$$

b) Examples:

i) The Sample Mean or Average

(a) The mean, i.e. μ or $E(X)$, is a population feature, the sample mean or average is an estimate of it.

(b) Def. $\bar{X} \equiv g(x_1, x_2, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i$

(c) Other measures of the sample's central tendency, the sample median, mode, midrange, etc. can be defined analogously to their population cousins as well.

ii) The Sample Standard Deviation:

(a) usually, defined $s_x \equiv \left[\frac{1}{n-1} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \right]^{1/2}$

• occasionally, we use (for reasons to be explained)

$\hat{\sigma}_x = \left[\frac{1}{n} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \right]^{1/2}$ instead

or divide by some other quantity (related to n) for other reasons

(b) other measures of the sample spread, sample mean absolute deviation, or sample range etc. can again be defined similarly to their population analogs.
Sample variance, e.g., is square of sample std. dev.

(iii) Sample Relationships

(a) sample covariance:

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

(b) sample correlation:

$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

(c) it is often useful to organize & arrange sample variances & covariances / or correlations in matrices:

(1) Sample Covariance Matrix

$$S_x = \begin{bmatrix} s_{x_1 x_1} & \dots & s_{x_1 x_n} \\ \vdots & \ddots & \vdots \\ s_{x_n x_1} & \dots & s_{x_n x_n} \end{bmatrix}$$

(2) Sample Correlation Matrix

$$R_x = \begin{bmatrix} r_{x_1 x_1} & \dots & r_{x_1 x_n} \\ \vdots & \ddots & \vdots \\ r_{x_n x_1} & \dots & r_{x_n x_n} \end{bmatrix}$$

- n.b.:
- Diagonal elements of S_x are variances of x_i , right?
 - Diagonal elements of R_x are all one (1), right?

III.B.2.c. Some useful short-cuts for calculating sample statistics.

(4)

$$i) S_x^2 \equiv \text{sample variance} = \frac{1}{n-1} \left[\left(\sum_{i=1}^n X_i^2 \right) - n \bar{X}^2 \right] = \sum \frac{1}{n-1} X_i^2 - \frac{n}{n-1} \bar{X}^2$$

(d. $S_x = \sqrt{S_x^2}$)

$$ii) S_{xy} \equiv \text{sample cov. of } x \text{ \& } y = \frac{1}{n-1} \left[\left(\sum_{i=1}^n X_i Y_i \right) - n \bar{X} \bar{Y} \right] = \sum \frac{1}{n-1} X_i Y_i - \frac{n}{n-1} \bar{X} \bar{Y}$$

$$iii) r(a \cdot x, b \cdot y) = \frac{a \cdot b}{|a \cdot b|} r_{x,y} \quad [e.g. r_{x,y} = .5, \Rightarrow r(2 \cdot x, -1 \cdot y) = -.5]$$

↑ constants ↑ absolute value

$$iv) s(a + b \cdot x, c + d \cdot y) = b \cdot d \cdot S_{xy} \quad [e.g. s(2 + 3x, 4 + 2y) = 12 \cdot S_{xy}]$$

↑ constants ↓ abs. val.

• implies, for example, $S^2(a + bx) = b^2 \cdot S_x^2 \Rightarrow s(a + bx) = |b| S_x$

• in short, all the useful properties of expectations, variances, & covariances we used for population P.V.'s last week work exactly analogously for their sample statistic analogs.

III. C. Sampling Distributions (of Sample Statistics)

- 1) The random sample itself can be seen as presenting a "sample distribution". This is not to be confused with the sampling distribution of a statistic.
- 2) Since each observation is a draw from a probability distribution, if we drew another sample, the observations would be different. Thus, across repeated samples, sample statistics calculated from those samples will vary in some way. The sampling distribution of a statistic is the p.f. or p.d.f. which describes how the statistic varies across these (usually hypothetical) samples.
- 3) Thus, we can speak, for example, of $\bar{X} \sim f(\bar{X}, \theta)$, i.e. of the sample average being distributed, across repeated samples, according to $f(\bar{X}, \theta)$. We might, furthermore, be interested in the expected value or variance, etc., of \bar{X} or whatever other statistic across these repeated samples.

3) Thus, we can speak, for example, of $\bar{X} \sim f(\bar{x}, \theta)$, i.e. of the sample average being distributed, across repeated samples, according to $f(\bar{x}, \theta)$. We might, furthermore, be interested in the expected value or variance, etc., of \bar{X} or whatever other statistic across these repeated samples.

Eq. $E(\bar{X}) = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right]$ by definition of \bar{X}

$= \frac{1}{n} E(\sum X_i)$ b/c can pull constants out of $E(\cdot)$

$= \frac{1}{n} (E(X_1) + E(X_2) + E(X_3) + \dots)$ writing out the sum

$= \frac{1}{n} (\mu + \mu + \mu + \dots + \mu)$ each X_i is a S.R.S. from population with some mean μ

$= \frac{1}{n} (n \cdot \mu)$

$= \mu$

$V(\bar{X}) = V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \left(\frac{1}{n}\right)^2 V\left(\sum_{i=1}^n X_i\right) = \left(\frac{1}{n}\right)^2 \sum_{i=1}^n V(X_i)$

$= \frac{1}{n^2} (n \cdot \sigma_x^2)$

$= \frac{\sigma_x^2}{n}$

nb. in this step, since X_i independent, $Cov(X_i, X_j) = 0$

III.C.3. Mean & Variance of the Sampling Distribution of \bar{X}

$$E(\bar{X}) = \mu$$

$$V(\bar{X}) = \frac{1}{n} \cdot \sigma_x^2$$

Notes: a) This is true for X 's drawn ^(independently) from any p.f. or p.d.f. with defined means & variances.

b) Variance of averages across repeated samples is $\frac{1}{n}$ (the variance of the underlying random variable)

\Rightarrow "averaging reduces variance" in this sense & by a factor of $\frac{1}{n}$ where n is the sample size being averaged.

4. Example: if $X_i \sim N(\mu, \sigma_x^2) \forall i$ (\Rightarrow the X_i are indep.)

$$\bar{X} \sim N(E(\bar{X}), V(\bar{X})) \quad \text{b/c sum of normals is normal}$$

$$\bar{X} \sim N\left(\mu, \frac{\sigma_x^2}{n}\right)$$

4. We will often be interested in how a statistic behaves across repeated samples as those (hypothetically repeated) samples get larger & larger:

e.g. a) $\lim_{n \rightarrow \infty} V(\bar{X}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$

b) $\lim_{n \rightarrow \infty} E(\bar{X}) = \lim_{n \rightarrow \infty} \mu = \mu$

c) so, as $n \rightarrow \infty$, the distribution of \bar{X} goes to a "spike" over μ .
more on this in a bit.

IV. (Point) Estimation of Parameters

A. First some Definitions

1. point estimate a statistic that gives a single-value estimate or "guess" for a parameter, θ , or for each of a set of parameters, Θ
2. standard error the standard deviation of the sampling distribution of some point estimate
 - the sampling standard deviation is the standard error; the square of that is the sampling variance
3. interval estimate a statistic that gives some range of values as a guess of θ or Θ . Usually, the interval is constructed in an appropriate (i.e. useful) manner from a point estimate and its standard error.

e.g. "margin of error"; 40% favor Clinton ± 2 : $\theta = \text{true \% favoring Clinton}$, 40% is the point estimate, $\hat{\theta}$ is the standard error, "margin of error" $\equiv 2 \cdot \hat{\theta}$. (so here, $\hat{\theta}$ must be 1) & the interval estimate is 40% ± 2 or $\{38\%, 42\%\}$

4. Estimator: an estimator is a rule or a strategy for guessing at a parameter; i.e. it is the function that produces an estimate:

e.g.: $g(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n X_i$ is an estimator for μ (that we often call \bar{x})

5. Finite (or "small") Sample: any fixed sample size less than ∞

Asymptotic (or "large" or "infinite" sample: a sample of infinite or "going toward infinite" size

IV.B. Criteria for Judging/Comparing Estimators

An estimator for (a) parameter(s), θ , can be any old function of the data, but surely some are "better" in some sense than others as "guesses" of θ . The "better" we are usually concerned with relates to how they perform across repeated samples; i.e., we compare them in terms of their sampling distributions.

1. Unbiasedness: An estimator of θ , call it $\hat{\theta}$, is unbiased iff $E(\hat{\theta}) = \theta$.

a. Examples: i. \bar{X} is an unbiased estimator of μ

$$E(\bar{X}) = E\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \frac{1}{n} E(\sum_{i=1}^n X_i) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} \cdot n \mu = \mu \quad \square$$

ii. $\hat{\theta} = 1^{\text{st}}$ observation of X, X_1 , is an unbiased estimator of μ

$$E(\hat{\theta}) = E(X_1) = \mu \quad \square$$

b. Bias($\hat{\theta}$) $\equiv E(\hat{\theta}) - \theta$ (= 0 in above two examples)

$\hat{\theta} = \frac{1}{3} X_1$ is a biased estimate of μ , its bias is

$$E\left(\frac{1}{3} X_1\right) - \mu = \frac{1}{3} E(X_1) - \mu = \frac{1}{3} \mu - \mu = -\frac{2}{3} \mu$$

2. Efficiency: An unbiased estimator, $\hat{\theta}_1$, is more efficient than another unbiased estimator, $\hat{\theta}_2$, iff $V(\hat{\theta}_1) < V(\hat{\theta}_2)$

e.g. $\hat{\theta}_1 \equiv \bar{X}$ is more efficient than $\hat{\theta}_2 \equiv X_1$: $V(\bar{X}) = \frac{\sigma_x^2}{n} < V(X_1) = \sigma_x^2$

Efficiency in the multivariate case:

$\hat{\theta}_1$ is more efficient than $\hat{\theta}_2 \Leftrightarrow V(\hat{\theta}_2) - V(\hat{\theta}_1)$ is positive definite;

$$\text{i.e., } \mathbf{x}' [V(\hat{\theta}_2) - V(\hat{\theta}_1)] \mathbf{x} > 0, \forall \mathbf{x}.$$

Note: efficiency is unambiguously desirable only when comparing two estimators with equal bias or consistency properties/performance. More generally, a tradeoff arises between bias/consistency and efficiency. A measure that combines these desiderata is:

3. Mean-Squared Error (MSE):

$$\text{MSE}(\hat{\theta}) \equiv E[(\hat{\theta} - \theta)^2] = \text{Var}(\hat{\theta}) + (\text{bias}(\hat{\theta}))^2$$

$$= E[\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2]$$

$$= E(\hat{\theta}^2) - 2E(\hat{\theta} \cdot \theta) + \theta^2$$

$$= [E(\hat{\theta})]^2 + V(\hat{\theta}) - 2E(\hat{\theta}) \cdot \theta + \theta^2$$

for any R.V., $E(X^2) = \mu^2 + \sigma^2$

for any 2 R.V.'s, X & Y ,
 $E(XY) = \text{Cov}(X, Y) + E(X)E(Y)$
 $\text{Cov}(\hat{\theta}, \theta) = 0$ b/c θ is a constant

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

$$\Rightarrow [\text{bias}(\hat{\theta})]^2 = [E(\hat{\theta})]^2 - 2E(\hat{\theta})\theta + \theta^2 \quad ; \text{ so } \textcircled{4} + \textcircled{5}$$

Two Ways to Show $E[(\hat{\theta} - \theta)^2] = V(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2$

$$\textcircled{1} E[(\hat{\theta} - \theta)^2] = V(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2$$

$$E(\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2) = E[(\hat{\theta} - E(\hat{\theta}))^2] + [E(\hat{\theta}) - \theta]^2$$

$$E(\hat{\theta}^2) - 2\theta^2 + \theta^2 = E(\hat{\theta}^2 - 2E(\hat{\theta})\hat{\theta} + [E(\hat{\theta})]^2) + E(\hat{\theta})^2 - 2E(\hat{\theta})\theta + \theta^2$$

$$\left[\begin{array}{l} \text{b/c for any RV} \\ \hat{\theta} = \theta + (\hat{\theta} - \theta) \\ \theta [\theta + (\hat{\theta} - \theta)] \\ \theta \quad + \theta \end{array} \right] \quad = \quad E(\hat{\theta}^2) - 2[E(\hat{\theta})]^2 + [E(\hat{\theta})]^2 + [E(\hat{\theta})]^2 - 2\theta^2 + \theta^2$$

taking $E(\theta)$ through

$$\left[\begin{array}{l} \text{b/c for any RV} \\ E(\hat{\theta}) = \theta + E(\hat{\theta} - \theta) \\ \theta \cdot [\theta + (E(\hat{\theta} - \theta))] \\ \theta^2 + 0 \end{array} \right]$$

$$E(\hat{\theta}^2) - \theta^2 = E(\hat{\theta}^2) - \theta^2$$

just add & subtract $E(\hat{\theta})$

$$\textcircled{1} E\{(\hat{\theta} - \theta)^2\} = E\{[(\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)]^2\}$$

$$= E\{(\hat{\theta} - E(\hat{\theta}))^2 + (E(\hat{\theta}) - \theta)^2 + 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)\}$$

$$= V(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2 + 2 \cdot E(\hat{\theta} \cdot E(\hat{\theta}) - \hat{\theta} \cdot \theta - E(\hat{\theta})^2 + E(\hat{\theta})\theta) \\ + 2 \cdot [(E(\hat{\theta}))^2 - E(\hat{\theta})\theta - [E(\hat{\theta})]^2 + E(\hat{\theta})\theta] \\ + 2 \times 0$$

$$= V(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2$$

(Minimum) Mean Squared-Error (“(best) combined min-bias & max-ffic.”):

$$E \left[(\hat{\theta} - \theta)^2 \right] = V(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2$$

4. More generally, we can assume / or it is appropriate, to apply some other function reflecting the costs of “mistakes”.

E.g. Min MSE assumes/applies costs of mistakes according to the function $(\hat{\theta} - \theta)^2$. We could imagine scenarios where costs of mistakes are zero for small mistakes & constant for any large error e.g.

$$L(\hat{\theta}) = \begin{cases} 0 & \forall \hat{\theta} - \theta < c \\ c & \forall \hat{\theta} - \theta \geq c \end{cases}$$

or any-old loss function we like

5. Minimum Variance Unbiased Estimator: one which is unbiased & has the lowest variance of any estimator that is likewise unbiased.

(MVUE)

Minimum Variance Linear Unbiased Estimator or “Best” Linear Unbiased Estimator:

(MVLUE)

(BLUE)

one which is unbiased & linear & has the lowest variance of any estimator that is likewise linear & unbiased.

• in general, we cannot know if some estimator is MVUE or MVLUE / BLUE. Similarly, we do not generally know if some estimator is Minimum MSE. However, Cramér & Rao proved an interesting result:

a) Proving **minimum-variance** (*most-efficient*) among some class of estimators (among unbiased estimators, or among linear & unbiased estimators...) is hard. One way can know you have min-var (or **best**) among unbiased is via **Cramér-Rao Lower-Bound**

Cramér-Rao Lower Bound \rightarrow for any "well-behaved" p.d.f., any unbiased estimator of parameter θ , call it $\hat{\theta}$, will have variance at least as large as:

$$[I(\theta)]^{-1} \equiv \left(-E \left[\frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \right] \right)^{-1}$$

$L(\theta)$ here is the "likelihood function" of the data, i.e., loosely, it is

$L(\theta) = \prod_{i=1}^n f(x_i, \theta)$; e.g. suppose x_i drawn from exponential distribution,

$x_i \sim f(x_i, \theta) = \theta e^{-\theta x_i}$ $\theta =$ hazard rate
 $\xrightarrow{\text{life-expectancy}} \frac{1}{\theta} = \text{mean}$
 $\frac{1}{\theta^2} = \text{var}$

then, $L(\theta) = \prod_{i=1}^n f(x_i, \theta) = \prod_{i=1}^n \theta e^{-\theta x_i} = \theta^n e^{-\theta \sum x_i}$

thus, $\ln L(\theta) = n \ln \theta - \theta \sum_{i=1}^n x_i$ $\hat{\theta}_{ML} = \left(\frac{\sum x_i}{n} \right)^{-1}$

$\therefore \frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \equiv \frac{\partial \left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)}{\partial \theta} = \frac{\partial \left(\frac{n}{\theta} - \sum x_i \right)}{\partial \theta} = -\frac{n}{\theta^2}$

So, the C-R lower bound = $\left(-E \left(-\frac{n}{\theta^2} \right) \right)^{-1} = \left[E \left(\frac{n}{\theta^2} \right) \right]^{-1} = \left[\frac{n}{\theta^2} \right]^{-1}$
 $= \frac{\theta^2}{n}$

Thus, C-R tells us that any ^{unbiased} estimator of the parameter θ in an exponential distribution must have a variance (across repeated samples) of at least $\left[\frac{\theta^2}{n} \right]$. That implies, for example, that if we can find an ^{unbiased} estimator with that variance, we can stop looking: no other unbiased estimator will do any better.

6. Large-Sample Criteria for Estimators:

a. Convergence in Probability

i. Def X_n converges in probability to a constant, c , written $X_n \xrightarrow{p} c$,
 $\Leftrightarrow \lim_{n \rightarrow \infty} \text{Prob}(|X_n - c| > \epsilon) = 0$ for any positive ϵ .
obs. value

ii. Def X_n converges in ^{something} quadratic mean (assuming X_n has mean μ & var σ_n^2)
 $\Leftrightarrow \lim_{n \rightarrow \infty} \mu = c$ & $\lim_{n \rightarrow \infty} \sigma_n^2 = 0$ & so $\text{plim } X_n = c$
(mean square) "pec. lim"

iii. 1 & 2 are not quite the same (2 is stricter), but for our purposes don't sweat the difference

IV. B. 6. b. Consistency: an estimator, $\hat{\theta}$, of a parameter, θ , is consistent

$$\Leftrightarrow \text{plim } \hat{\theta} = \theta$$

• The point: a consistent estimator goes to exactly the right answer (i.e. exactly the parameter you are estimating) as the sample size goes to infinity

• Example: \bar{X} is a consistent estimator of μ

$$\left. \begin{array}{l} \textcircled{1} E(\bar{X}) = \mu \\ \textcircled{2} V(\bar{X}) = \sigma^2/n \end{array} \right\} \Rightarrow \begin{array}{l} \textcircled{1} \lim_{n \rightarrow \infty} E(\bar{X}) = \mu \\ \textcircled{2} \lim_{n \rightarrow \infty} V(\bar{X}) = 0 \end{array}$$

$\therefore \bar{X}$ converges in quadratic mean to μ . (which is stricter than and therefore) $\Rightarrow \text{plim } \bar{X} = \mu$

• Extensions: ① $p \lim \frac{1}{n} \sum_i g(x_i) = E(g(x))$

e.g. $E(X^2) = \Theta$
 $\Theta = \frac{1}{n} \sum_i x_i^2$

$\Rightarrow \hat{\Theta}$ is consistent

i.e. the average of any function evaluated at each observation is a consistent estimator of the mean of that function.

② (Slutsky Theorem) for any continuous $g(x_n)$, $p \lim g(x_n) = g(p \lim x_n)$

examples: $p \lim 3x^2 = 3(p \lim X)^2$

$p \lim \frac{\bar{X}^2}{s^2} = \frac{(p \lim \bar{X})^2}{(p \lim s)^2} = \frac{\mu^2}{\sigma^2}$

C. Convergence in Distribution & Limiting Distribution

i. Def X_n converges in distribution to a random variable X with cumulative probability function, $F(x)$ \iff (different)

$\lim_{n \rightarrow \infty} |F_n(x) - F(x)| = 0$ (written $X_n \xrightarrow{d} F(x)$ or $X_n \xrightarrow{d} X$)
 ↑ abs. value.

i.e. as $n \rightarrow \infty$, the distribution of X_n becomes ever more like, & in the limit exactly equal to, some (other) distribution

• The distribution being "approached" in this way is called the limiting or asymptotic distribution of X_n .

ii. The most common & most important example: $X_n \stackrel{d}{\rightarrow} N(0,1)$

iii. Useful Combinations of the Two Convergences:

$$\begin{aligned} & \textcircled{1} X_n \xrightarrow{d} X \text{ \& } \text{plim } Y_n = C \\ & \textcircled{2} X_n \xrightarrow{d} X \text{ \& } g(X_n) \text{ continuous} \\ & \Rightarrow g(X_n) \xrightarrow{d} g(X) \end{aligned} \Rightarrow \begin{aligned} X_n Y_n & \xrightarrow{d} CX \\ X_n + Y_n & \xrightarrow{d} X + C \\ X_n / Y_n & \xrightarrow{d} X/C \quad (C \neq 0) \end{aligned}$$

iv. Another Useful Example: $F_{1,n} \xrightarrow{d} \chi^2_1$

v. Convergence in Distribution is also sometimes written

$$X \sim f(x) \quad \text{"X is asymptotically distributed f(x)"}$$

IV. B. 6. d. Best (ie. Minimum Variance) Asymptotic Normal Estimator Consistent (BANC)

$$\hat{\theta} \xrightarrow{d} N(\theta, \phi^2/n)$$

$\phi^2 \leq \phi^{*2}$ for any alternative estimator which converges in distribution to $N(\theta, \phi^{*2}/n)$

• The Reason we care about the CLT & limiting distributions is that we (more or less, speaking a bit loosely) "use" them to talk about the asymptotic distributions of statistics we care about.

e.g., as already noted, $\bar{X}_n \overset{\Delta}{\sim} N(\mu, \sigma^2/n)$
& $Z = ((\bar{X}_n - \mu) / (\sigma/\sqrt{n})) \overset{\Delta}{\sim} N(0, 1)$

• Z is not exactly distributed $N(0, 1)$, but

- 1) As n (sample size gets larger) it is more & more like it
- 2) The closer X (i.e. the actual, exact distribution of X) is to normal to begin with, the sooner & quicker the distribution of Z goes to $N(0, 1)$

Consider, for example, the statistic: $(\hat{\beta} - \beta_0) / \text{s.e.}(\hat{\beta}) = t$

i.e., the stat we use to test the hypothesis that $\hat{\beta} = \beta_0$. The CLT & Asymptotic Distribution theory tell us that if

- a) $\hat{\beta}$ is an unbiased estimator of β & if
- b) $\text{s.e.}(\hat{\beta})$ is a consistent estimator of $\sigma_{\hat{\beta}}$, then

$$t \overset{\Delta}{\sim} N(0, 1)$$

• Now, if, as in fact we often will, we know a bit more about $\hat{\beta}$ & $\text{s.e.}(\hat{\beta})$ we might know their actual distribution (e.g., under OLS assumptions, $\epsilon \sim \tau_n$), but we don't actually need anything more than "good" estimates of $\hat{\beta}$ & $\text{s.e.}(\hat{\beta})$ + the C.L.T. to get "asymptotic normality" for the "t" statistic.

B. The Law of Large Numbers (LLN) & the Central Limit Theorem (CLT) (the two most-important theorems in probability & statistics):

The Law of Large Numbers: simple random sampling from any population with $E(X) = \mu$ & $V(X) = \sigma^2$, the average, $\bar{X} \xrightarrow{p} \mu$

example: $X =$ R.V. equal to a die roll's upward face;

$$\bar{X}_n \xrightarrow{p} 3.5 \quad (\text{i.e. } \lim_{n \rightarrow \infty} \frac{1}{n} \sum (\text{die rolls}) = E(\text{die roll}))$$

↑
number of die rolls

Central Limit Theorem: in S.R.S. from any population w/ mean μ & variance σ^2 ,

$$\left\{ Z = \frac{\bar{X}_n - \mu}{(\sigma/\sqrt{n})} \right\} \xrightarrow{d} N(0, 1)$$

1. This is the most useful form of the CLT, notice also, though, that the left-hand side can be rewritten:

$$\sqrt{n} \cdot \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \xrightarrow{d} N(0, 1)$$

2. Given our rule that $V(bX) = b^2 V(X)$, we can also rewrite it:

$$\sqrt{n} \cdot (\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2) \quad (\text{as Greene Does})$$

3. Given our rule that $E(b+X) = b + E(X)$, we can also rewrite it

$$\bar{X}_n \xrightarrow{d} N\left(\mu, \frac{\sigma^2}{n}\right)$$

if X is drawn from population w/ mean μ , but variance which is different for every observation, $\sigma_{X_i}^2$, then

$$\sqrt{n} \cdot (\bar{X}_n - \mu) \xrightarrow{d} N(0, \bar{\sigma}^2)$$

$\bar{\sigma}^2$ is the "average variance," $\frac{1}{n} \sum \sigma_i^2$

Multivariate Extensions of the CLT:

1. if X_1, \dots, X_n are a S.R.S. from a joint distribution with mean μ and variance Q , then $\sqrt{n} (\bar{X}_n - \mu) \xrightarrow{d} N(0, Q)$

matrix (positive definite)

& all the permutations analogous to the above.

2. if X_1, \dots, X_n S.R.S. with mean μ , & variance $(X_i) = Q_i$

(differs for each i)

$$\sqrt{n} (\bar{X}_n - \mu) \xrightarrow{d} N(0, Q)$$

where $\bar{Q}_n = \frac{1}{n} \sum Q_i$

$$\text{d } Q = \lim_{n \rightarrow \infty} \bar{Q}_n$$

3. That is, the CLT generalizes in a more or less transparent way to the multivariate case

III. Basic Statistics Review, Summary:

- A. A *statistic*, or an *estimate*, is any quantity calculated from the data.
- B. An *estimator* is formula or strategy applied to calculate an *estimate*.
- C. Properties of estimators:

1. **Unbiasedness** (“right on average, across samples”): $E(\hat{\boldsymbol{\theta}}) = \boldsymbol{\theta}$.

2. **Consistency** (“exactly right, w/o variance, as sample-size $\rightarrow\infty$ ”): $\lim_{n\rightarrow\infty} \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$.

3. **Efficiency** (“all useful information used; no wasted info”):

$\hat{\boldsymbol{\theta}}_1$ is more efficient than $\hat{\boldsymbol{\theta}}_2 \Leftrightarrow V(\hat{\boldsymbol{\theta}}_2) - V(\hat{\boldsymbol{\theta}}_1)$ is positive definite;

$$\text{i.e., } \mathbf{x}' [V(\hat{\boldsymbol{\theta}}_2) - V(\hat{\boldsymbol{\theta}}_1)] \mathbf{x} > 0, \forall \mathbf{x}.$$

4. **(Minimum) Mean Squared-Error** (“(best) combined min-bias & max-effic.”):

$$E\left[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^2\right] = V(\hat{\boldsymbol{\theta}}) + [bias(\hat{\boldsymbol{\theta}})]^2$$

5. **Minimum-Variance (Best) Unbiased Estimator** (BUE) & **Minimum-Variance (Best) Linear Unbiased Estimator** (BLUE), **Minimum-Variance (Best) Asymptotically Normal & Consistent Estimator** (BANC): Self-explanatory...

6. Following can exemplify *estimator, estimate, unbiased, consistent, efficient*:

E.g.

$$\begin{aligned}
 E(\bar{X}) &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] && \text{by definition of } \bar{X} \\
 &= \frac{1}{n} E(\sum X_i) && \text{b/c can pull constants out of } E(\cdot) \\
 &= \frac{1}{n} (E(X_1) + E(X_2) + E(X_3) + \dots) && \text{writing out the sum} \\
 &= \frac{1}{n} (\mu + \mu + \mu + \dots + \mu) && \text{each } X_i \text{ is a S.R.S.} \\
 &= \frac{1}{n} (n \cdot \mu) && \text{from population with} \\
 &= \mu && \text{same mean } \mu
 \end{aligned}$$

$$\begin{aligned}
 V(\bar{X}) &= V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \left(\frac{1}{n}\right)^2 V\left(\sum_{i=1}^n X_i\right) = \left(\frac{1}{n}\right)^2 \sum_{i=1}^n V(X_i) \\
 &= \frac{1}{n^2} (n \cdot \sigma_x^2) \\
 &= \frac{\sigma_x^2}{n}
 \end{aligned}$$

n.b. in this step, since X_i independent, $\text{Cov}(X_i, X_j) = 0$

IV. Maximum-Likelihood Estimation:

A. We now actually have more or less all the tools we need for Maximum Likelihood Estimation (MLE)

1. We view outcomes as being produced by some p.f. or p.d.f., $f(x_i, \theta)$
2. If each outcome is i.i.d. (i.e., is from a S.R.S.), then their joint distribution is the product of their marginal distributions:

$$f(x_1, x_2, x_3, \dots, x_n, \theta) = f(x_1, \theta) f(x_2, \theta) f(x_3, \theta) \dots f(x_n, \theta)$$

3. What we do now is choose θ (i.e. estimate it) so that it maximizes the "probability" or likelihood of our having observed the data we actually have observed.

B. That is, given the data, X , define the (joint) conditional distribution of θ given X . Maximize this with respect to θ .

C. Examples:

1. Each $X_i \sim \text{Poisson}(\theta)$: $f(x_i, \theta) = \frac{e^{-\theta} \theta^{x_i}}{x_i!}$

(A Poisson Distribution describes a count (R.V.) of events happening at average rate θ per fixed period)

a. $f(x_1, \dots, x_n, \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!}$

b. $\text{Max}_{\theta} f(x, \theta) = \text{Max}_{\theta} \ln f(x, \theta) = \text{Max}_{\theta} \sum_{i=1}^n [\ln(e^{-\theta} \theta^{x_i}) - \ln(x_i!)]$

$= \text{Max}_{\theta} \sum_{i=1}^n (-\theta + x_i \ln \theta - \ln(x_i!))$
call this $L(\theta)$

First-order Condition $\frac{\partial L(\theta)}{\partial \theta} = 0$

Aside: $\text{Max} f(x) = \text{Max} g(f(x))$

for any $g(\cdot)$ which is a strictly monotonic increasing function such as $g(\cdot) = \ln(\cdot)$

Aside #2: "strictly monotonic increasing" means that, in $g(x)$, if x goes up so does $g(x)$

VI. c. b. d. Maximizing Log-Likelihood of Poisson Model

$$L(\theta) = \sum_{i=1}^n (-\theta + x_i \ln \theta - \ln x_i!) = -n\theta + \left(\sum_{i=1}^n x_i\right) \ln \theta - \sum_{i=1}^n \ln(x_i!)$$

$$\text{F.O.C. } \frac{\partial L(\theta)}{\partial \theta} = 0 \Rightarrow -n + \frac{\sum_{i=1}^n x_i}{\theta^*} = 0 \Rightarrow \theta^* = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}$$

e. Second-Order Condition:

$$0 > \frac{\partial \left(\frac{\partial L(\theta)}{\partial \theta} \right)}{\partial \theta} \Rightarrow \frac{\partial (-n + \theta^{-1} \cdot \sum x_i)}{\partial \theta} = -\theta^{-2} (\sum x_i) < 0 \quad \checkmark$$

↳ greater than zero

f. Thus, MLE for Poisson parameter, θ , is $\hat{\theta} = \bar{X}$.

VI. C. 2. MLE Example 2: MLE for Normal

$$a. f_n(x, \theta) = (2\pi\sigma^2)^{-1/2} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$b. L(\theta) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\ln L(\theta) = \sum_{i=1}^n \ln \left[(2\pi\sigma^2)^{-1/2} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right]$$

$$= \sum_{i=1}^n \left[-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \sum_{i=1}^n \left(-\frac{1}{2} \ln 2 - \frac{1}{2} \ln \pi - \frac{1}{2} \ln \sigma^2 - \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$= -\frac{n}{2} \ln 2 - \frac{n}{2} \ln \pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 (\sigma^2)^{-1}$$

$$c. \text{F.O.C. } \frac{\partial \ln L(\theta)}{\partial \theta} = 0:$$

$$i.) \frac{\partial \ln L(\theta)}{\partial \mu} = (\sigma^2)^{-1} \sum_{i=1}^n (x_i - \mu) = 0$$

$$ii.) \frac{\partial \ln L(\theta)}{\partial \sigma^2} = \frac{n}{2\sigma^2} + \frac{1}{2} \cdot (\sigma^2)^{-2} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

• Multiply both sides of (i) by $(\sigma^2)^{-1} \Rightarrow \sum_{i=1}^n (x_i - \mu) = 0$
 $\Rightarrow \sum_{i=1}^n x_i - n\mu^* = 0$
 $\Rightarrow \mu^* = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$

• Substitute that into (ii) and you'll find $\sigma^{*2} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

VI. C. 3. We note (in passing for now) that Maximum Likelihood Estimation is especially useful when we have some substantive model for the parameter(s), such as $\mu = a + bX$ or something. It's especially especially useful when the probability function underlying the outcomes we want to model is non-linear.

VI. D. Properties of Max Like Estimators:

- 1) Consistent ($\text{plim } \hat{\theta}_{ML} = \theta$)
- 2) Asymptotically Normal ($\hat{\theta}_{ML} \overset{\Delta}{\sim} N(\theta, \{I(\theta)\}^{-1})$)
- 3) Asymptotically Achieves Cramér-Rao Lower Bound for Consistent Estimators (i.e. is "best")
- 4) Invariant: if $\hat{\theta}_{ML}$ is MLE for θ , then $g(\hat{\theta}_{ML})$ is MLE for $g(\theta)$
 - ↳ eg. in normal example above, MLE for σ is square root of MLE for σ^2

VII. Interval Estimation & Hypothesis Testing

A. Some Preliminary Facts We Use a lot in Interval Est. & Hypoth. Testing:

Note: the process for hypoth. testing is to determine the distribution of a test statistic under the null hypothesis & see how "odd" the actual stat. is in the case at hand. For intervals, it's to see how wide an interval must be to be reasonably

1. As we know, in sampling X from a normal, with known mean μ & variance σ^2 ,

$$\frac{X - \mu}{\sigma} \sim N(0, 1)$$

• this is exact, not an approximate distribution

sure to contain the true parameter
• note on this in a bit

2. When taking a sample average from a known $N(\mu, \sigma^2)$:

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

• also exact

3. When sampling X from an unknown distribution with known variance,

$$\frac{X - \mu}{\sigma} \xrightarrow{d} N(0, 1)$$

• not exact, but converges to Normal

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

4. When sampling from a normal distribution with unknown variance:

$$\frac{X - \mu}{s} \sim t_{n-1} \xrightarrow{d} N(0, 1)$$

$$\frac{(\bar{X} - \mu)/(s/\sqrt{n})}{1} \sim t_{n-1} \xrightarrow{d} N(0, 1)$$

5. When sampling from an unknown distribution with unknown parameters

$$\frac{X - \mu}{s} \xrightarrow{d} N(0, 1)$$

$$\frac{\bar{X} - \mu}{s/\sqrt{n}} \xrightarrow{d} N(0, 1)$$

upshot: $\frac{X - \mu}{\sigma/\sqrt{n}}$ always $N(0, 1)$ in the limit (by CLT)
Only a question of how large sample must be for $N(0, 1)$ to be reasonable approx. Under some circumstances we know the $N(0, 1)$ happens along

$$6. V = \sum_{i=1}^n X_i^2 \quad \text{when each } X_i \sim N(0, 1) \Rightarrow V \sim \chi_n^2$$

• if (each $X_i \sim N(0, 1)$) then $V \sim \chi_n^2$

$$7. W = \left\{ \frac{(V_1/n_1)}{(V_2/n_2)} \right\} \quad \text{where } V_1 \sim \chi_{n_1}^2 \text{ and } V_2 \sim \chi_{n_2}^2 \Rightarrow W \sim F_{n_1, n_2}$$

• if each $V \sim \chi_n^2$, then $W \sim F_{n_1, n_2}$

$$\text{e.g. : } R^2 = 1 - \frac{(\sum e_i^2)}{(\sum y_i^2)}$$

under certain conditions, $R^2 \sim F$

B. Constructing Confidence Intervals:

1. We start with some estimate (statistic) which we "standardize" or "normalize" (i.e. subtract mean & divide by standard deviation):

$$\text{estimate} \rightarrow \bar{x} \quad \left\{ \frac{(\bar{x} - \mu)}{(s/\sqrt{n})} \right\} \leftarrow \text{standardization}$$

std dev. of \bar{x} (estimated)

2. Then we ask what "bounds" would contain this with some probability, given its distribution (here t_{n-1} by (4) above); i.e., we look for z 's that satisfy $\text{Prob}(-z < \frac{\bar{x} - \mu}{s/\sqrt{n}} < z) = (1 - \alpha)$

$$2. \quad \text{Prob} \left(-z < \left\{ \frac{\bar{x} - \mu}{s/\sqrt{n}} \right\} < z \right) = (1 - \alpha)$$

\rightarrow confidence level, e.g. $\alpha = .10$
 \Rightarrow 90% confidence interval for $\{ \cdot \}$ is $(-z, z)$

3. More usually, we are interested in an interval anchored on our estimate, \bar{x} here, & going +/- on either side of that so that the interval might encompass the parameter we're trying to estimate.

So we can rearrange (2.):

$$\text{Prob} \left(\underbrace{\bar{x} - \frac{z \cdot s}{\sqrt{n}}}_{\text{lower bound}} < \underbrace{\mu}_{\text{parameter}} < \underbrace{\bar{x} + \frac{z \cdot s}{\sqrt{n}}}_{\text{upper bound}} \right) = \underbrace{(1 - \alpha)}_{\text{confidence level}}$$

Steps of that rearrangement:

$$\text{Pr} \left(-z < \frac{\bar{x} - \mu}{s/\sqrt{n}} < z \right) = \text{Pr} \left(-\frac{z \cdot s}{\sqrt{n}} < \bar{x} - \mu < \frac{z \cdot s}{\sqrt{n}} \right) = \text{Pr} \left(-\frac{z \cdot s}{\sqrt{n}} - \bar{x} < -\mu < \frac{z \cdot s}{\sqrt{n}} - \bar{x} \right) = \text{Pr} \left(\bar{x} - \frac{z \cdot s}{\sqrt{n}} < \mu < \bar{x} + \frac{z \cdot s}{\sqrt{n}} \right)$$

4. Notes:

- a. We're implicitly assuming that we won't a symmetric confidence level ~~level~~ interval
- i.) if distribution of the statistic is symmetric, then the symmetric c.i. is also the smallest c.i. from lower to upper bound
 - ii.) if stat. is not symmetrically distributed, may be reasonable to get the smallest bound-range instead: this is complicated (I rarely done because of that, but we may try it later)
- b. Notice that it is the bounds that vary across repeated samples in 3, not the parameter. (The "truth" doesn't vary.) Thus the confidence interval tells us that in $x\%$ of repeated samples, a C.i. so constructed would contain the true parameter. (Tortuously close to what we'd like to say, but not, huh? In Bayesian analysis, FYI, you can say exactly what you want to say, but that's beyond us right now.)

5. Mechanics (Example)

\bar{X} taken from $x_i \sim \text{i.i.d. } N(\mu, \sigma^2)$ $i = \{1, \dots, 10\}$ ↖ unknown

$\bar{X} = .02$ $s = .1$ from our sample estimates

$$\Pr\left(\bar{X} - \frac{z \cdot s}{\sqrt{n}} < \mu < \bar{X} + \frac{z \cdot s}{\sqrt{n}}\right) = (1 - \alpha)$$

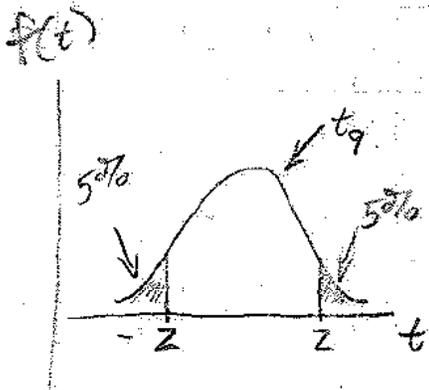
$$\Pr\left(.02 - z \frac{.1}{\sqrt{10}} < \mu < .02 + z \frac{.1}{\sqrt{10}}\right) = (1 - \alpha)$$

Suppose we want 90% c.i. (i.e., $\alpha = .1$), where does z come from? What's distrib. of $(\bar{X} - \mu) / (s / \sqrt{n})$? t_{n-1} or t_9 here. t is symmetric, so we want 5% on each side
 $\Rightarrow z$ value of 1.833 from a t -table (or computer program)

$$\Rightarrow \left[.02 - 1.833 \cdot \frac{.1}{\sqrt{10}}, .02 + 1.833 \cdot \frac{.1}{\sqrt{10}} \right]$$

$$[-.038, .078] \text{ is the 90\% c.i.}$$

So, 90% of intervals likewise constructed across repeated samples would contain the true mean.



→ practical rule of thumb: when in doubt, use the t

6. Final Note: Often when the distribution of the sample statistic is asymptotically normal, we use the t distribution rather than jump right to using the normal as an approximation. The t is clearly more conservative than $N(0, 1)$, but we can't actually say this is warranted nonetheless, except in the limit (asymptotically).

2. Type I & Type II Errors, Size, Power, Bias & Consistency of Tests

- a. Type I Error: Rejection of Null Hypothesis when it is in fact true.
- b. Type II Error: Failure to reject the Null Hypothesis when it is in fact false.
(Type III Error: confusing a Type I & Type II error -- (a geek joke))

c. The probability of a Type I error is the size (α -level, p -level, significance level) of a test.

d. The probability that a test will correctly reject the null, i.e. $1 - \text{Prob}(\text{Type II Error})$, is the power of a test. call this β

e. In general, if we are using all available info. correctly, ^{in estimating,} and we have correctly constructed the test to make best use of this information & estimate, then there is no way to make Type I errors less likely without making Type II errors more likely, and vice versa.

Example: if the criminal justice system is efficiently using available evidence, there is no way to make conviction of criminals more likely without also making conviction of the innocent more likely & v.v., there is no way to make conviction of innocent less likely without also making conviction of guilty less likely.

Similarly, for fixed efficiency \Rightarrow barring changes in efficiency, Miranda Rights, for example, leads to more innocent going free and more guilty going free.

- f. A Test is Uniformly Most Powerful if it has higher power than any other test of the same size for any admissible value of the true Parameter. Such tests are rarely possible, even more rarely found, what we generally do is compare the power of specific alternative tests.
- g. A test that is always (i.e. for any permissible parameter value) more likely to reject the null when it is false than when it is true is called unbiased.
- h. A test is consistent when its power $\rightarrow 1$ as the sample size increases. I.e., as sample size $n \rightarrow \infty$, the test always rejects a false null.
- in general, if the estimate being used in the test is consistent, then the test is consistent. Why? b/c consistent estimate of $\theta \Rightarrow \lim_{n \rightarrow \infty} \hat{\theta} = \theta$, so in the limit $\hat{\theta}$ converges to the true parameter & we can see directly whether H_0 is true or not.
- ← ~~is~~ closely related to unbiased estimates

Additional notes:

f. "...for specific (ranges of) parameter values."

g. As with h. regarding consistency: essentially, if estimate unbiased, and verdict directly related to estimate, then test unbiased.

3. Mechanics of Conducting Hypothesis Tests:

a. $H_0 : \theta = \theta_0$ Example: $H_0 : \beta = 0$
 $H_1 : \theta \neq \theta_0$ $H_1 : \beta \neq 0$

b. Convert Estimate, $\hat{\theta}$, to standardized form as in conf. intervals, under null being true

$$T = \left\{ \frac{\hat{\theta} - \theta_0}{\text{std. err.}(\hat{\theta})} \right\} \sim^A N(0, 1)$$

Example: $T = \frac{\hat{\beta} - 0}{\text{se}(\hat{\beta})} \sim^A N(0, 1)$

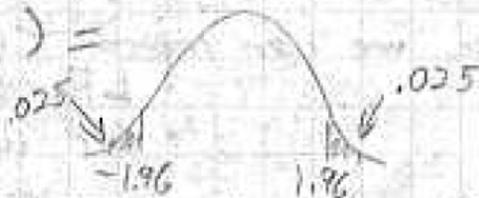
or $T = \frac{\hat{\beta}}{\text{se}(\hat{\beta})}$

as before, we usually use t_{n-1} to approximate rather than the asymptotic distribution $N(0, 1)$

c. Check that test statistic, t , against its distribution or approximate distribution, here t record the probability of T 's more extreme (greater in absolute value) than your actual T

Example $T = 1.96$, distribution t_{1000} ,

$$\text{Prob.}(|T| \geq 1.96) = .025 + .025 = .05$$



d. This is your p-level. It is the probability of having observed a test statistic, here T , this far or farther from the null, here 0, if the null were in fact true.

- If the p-level is \leq the predetermined size or α -level of the test, we reject the null (in favor of the alternative).

- note the 1-for-1 correspondence of hypoth. testing & confidence intervals: if reject, then c.i. did not contain null and vice versa.

e. Interpretational Note:

- The p-level does not tell us the probability the null is true (in non-Bayesian analysis). The null is either true (with probability 1) or not true (with probability 1). The p-level tells us the probability we would have observed a test stat. this far or farther from the null if the null were true. From this we infer that the null was likely untrue or not; the facts, however, are that the null either is or is not true (with $p=1$). n.b. use likely not probably

f. Rob's Rules: the standard procedure is to set the size of the test, α -level, to .10 or .05 or .01 & then stop anything that makes this or lower p-level. That's fine, but with modern tech, it's simple (the computer'll give it) to know the "exact" p-level. This is more info, we have the info, therefore under distributional assumptions/approximations give it.

g. The Broad (A Better by Rob's Rules) Interpretation:

- The question is "Is the estimated quantity, e.g. the relationship between X & Y, 'far' from the null, e.g. 0?"
- "Far" is defined in standard deviation units & how "significant" any increment of "far"ness is, i.e. how tortuous the straining of credulity given the evidence, i.e. how "long" a standard deviation is, is measured in increments of likelihood.
- Thus, p-levels tell us just how far we're straining credulity given the evidence of what has happened if we cling to the null. The lower the p-level, the harder is the stretch of imagination required to continue believing the null, to swallow.

VII. D. Three Common Types of Tests in Regression Analysis.

1. Think of hypotheses regarding (a) parameter(s), θ , such as

$$H_0: c(\theta) = 0$$

\hookrightarrow some function of theta

$$H_1: c(\theta) \neq 0$$

2. Likelihood-Ratio (LR) Tests:

a) if $c(\theta) = q$ ^{constants} is true, then imposing it to begin with (prior to estimation) should make little difference to the likelihood of the data (or fit of the regression)

b) \therefore , we ought to be able to base a test on $L(\hat{\theta}_N) / L_R(\hat{\theta}_N)$ $\leftarrow c(\theta) = q$ not imposed / imposed
or, equivalently, $\ln L(\hat{\theta}_N) - \ln L_R(\hat{\theta}_N)$ (just applying \ln rules)

c) In fact, under certain conditions we won't get into

$$\text{now } LR = 2 \cdot [\ln L(\hat{\theta}_N) - \ln L_R(\hat{\theta}_N)] \sim \chi_n^2$$

d) So, we can reject $c(\theta) = q$ if $\uparrow n = \text{number of conditions given by } c(\theta) = q$
 • LR larger than χ_n^2 for some predetermined α
 • or we can report the p-level of the test as before

"Degradation of Fit"
TESTS:
LR, DR²,
 ΔSSE

3. Wald Tests

- a) if $c(\hat{\theta}) = q$ true, then $c(\hat{\theta})$ should not be too "far" from q .
constant(s)
a "good" estimate, such as the MLE
- b) \therefore we ought to be able to find a test based on $c(\hat{\theta}) - q$ not being too far from zero.
- c) in fact, under certain conditions we won't go into now:

$$W = \{ (c(\hat{\theta}) - q)' (\text{Var}[c(\hat{\theta}) - q]) (c(\hat{\theta}) - q) \}$$

$$W \sim \chi_n^2$$

- d) reject or not, or report p-level from χ_n^2 exactly as before.
again, number of conditions given by $c(\hat{\theta}) = q$

4. Lagrange Multiplier Tests:

- a) if $c(\hat{\theta}) = q$ true, then imposing restraint should not have much effect in our optimizing the likelihood.
 In particular, the first-order condition that $\frac{\partial L(\theta)}{\partial \theta} = 0$ should not be "far" from holding.

- b) \therefore should be able to base a test on $\left\{ \frac{\partial \ln L(\theta)}{\partial \theta} \right\}$ with $c(\hat{\theta}) = q$ imposed being "near" zero.

- c) in fact, under conditions we're not going into, etc etc.

$$LM = \left[\frac{\partial \ln L_R(\hat{\theta}_R)}{\partial \hat{\theta}_R} \right]' [I(\hat{\theta}_R)]^{-1} \left[\frac{\partial \ln L_R(\hat{\theta}_R)}{\partial \hat{\theta}_R} \right]$$

Info matrix from Cramer Rao

$$LM \sim \chi_n^2 \text{ number of conditions again.}$$

And, as mentioned previously, can sometimes also conduct LM tests by clever auxiliary regression in which coefficients are the Lagrangian multipliers for the constraints.