

STATISTICS REVIEW

Background Information.

I. Measures of Return and Risk.

A. Operational definition of Ex Post, Nominal Return.

$$r = (\text{ending value} - \text{beginning value} + \text{cash flows}) / (\text{beginning value})$$

$$= (\text{ending} - \text{beginning}) / (\text{beginning}) + (\text{cash flows}) / (\text{beginning})$$

$$= (\% \text{ Capital Appreciation}) + (\% \text{ Dividend Income}).$$

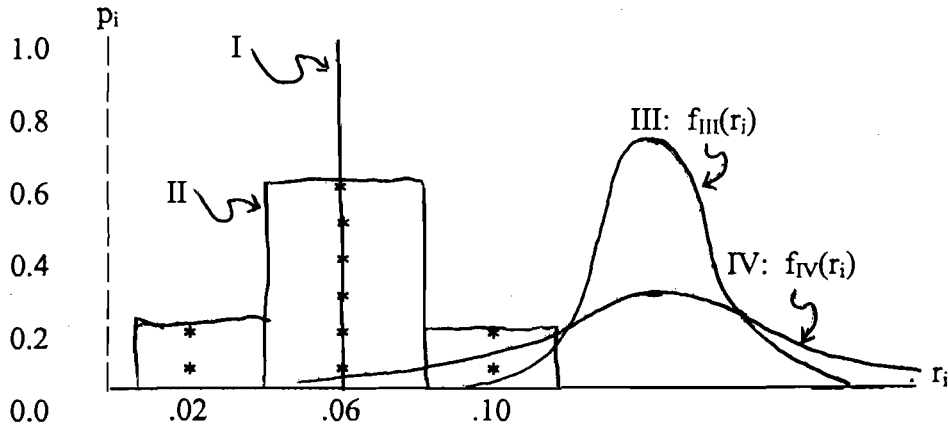
1. This is a measure of ex post, actual return earned in the past.
 - a. Accounting measure of past performance.
2. For investment decisions by individuals & management, need to consider expected future performance!
3. To get a measure of expected future return, must consider possible future outcomes.
4. This requires the probability distribution function (pdf) that reflects your expectations of possible future outcomes.
5. Consideration of this pdf introduces uncertainty – risk.

B. Expected Return and Risk.

– expected outcome & uncertainty about future return from investment.

1. Reflected in investor's probability distribution function (pdf).
2. Consider the pdf's of two possible investments, I & II:

Investment:	I	II	
Possible r (r _i)	.06	.02	.06 .10
Probability (p _i)	1.0	0.2	0.6 0.2



3. Note: I is a T.Bill; r = .06 with certainty (no risk).
II depends on future states of the world; uncertainty!
4. **Expected Return** = $E(r_i) = \sum p_i r_i$; i indexes states of world.
 - a. For investment I, $E(r_i) = (1.0)(.06) = .06$
 - b. For investment II, $E(r_i) = (.2)(.02) + (.6)(.06) + (.2)(.10) = .06$
5. Observe, expected return is the same for I & II, but II has more risk.
6. In reality, investor's pdf is continuous, anything between $-\infty$ & $+\infty$.
 - a. Thus, pdf is more like smooth curve such as III or IV above.
 - b. For continuous pdf's, expected return is defined as the integral:

$$E(r_i) = \int_{-\infty}^{\infty} f(r_i) r_i dr_i; \quad \text{analogous to } \sum p_i r_i.$$

C. How do we measure Risk? Several possibilities.

1. Range = $\text{Max}\{r_i\} - \text{Min}\{r_i\}$.

a. Problem: as sample size N increases, Range increases w/o bound.

2. Semi-Interquartile Range = $(X_{.75} - X_{.25}) / 2$
= $(75^{\text{th}} \text{ \% -tile} - 25^{\text{th}} \text{ \% -tile}) / 2$.

a. This does not suffer from the problem with the Range.

b. Used when the variance does not exist.

3. Variance = $\sigma^2 = E[r_i - E(r_i)]^2 = \sum p_i [r_i - E(r_i)]^2$.

a. For I, $\sigma^2 = (1.0)(.06 - .06)^2 = 0$. No uncertainty, no risk.

b. For II, $\sigma^2 = (.2)(.06 - .02)^2 + (.6)(.06 - .06)^2 + (.2)(.10 - .06)^2$
= $.00032 + 0 + .00032$
= $.00064$

c. Standard Deviation, σ .

For I, $\sigma = 0$; For II, $\sigma \approx .0253$

d. For continuous case, $\sigma^2 = \int_{-\infty}^{\infty} f(r_i) [r_i - E(r_i)]^2 dr_i$.

4. If r_i deviates further from its mean, distribution is more spread out:

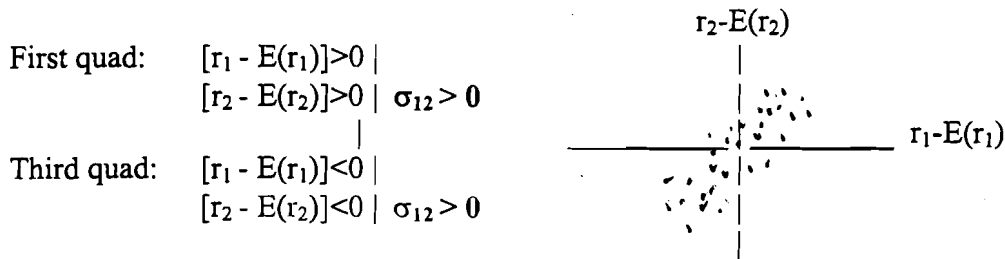
$[r_i - E(r_i)]$ is larger; $[r_i - E(r_i)]^2$ is larger; σ^2 is larger.

D. Covariance.

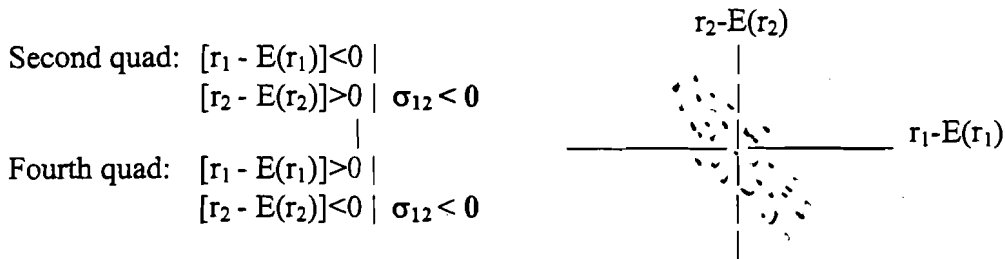
1. Defn: $\sigma_{12} = \text{Cov}(r_1, r_2) = E[r_1 - E(r_1)][r_2 - E(r_2)]$

2. Operational Defn: $\sigma_{12} = \sum_{i=1}^n [r_{1i} - E(r_1)][r_{2i} - E(r_2)]$

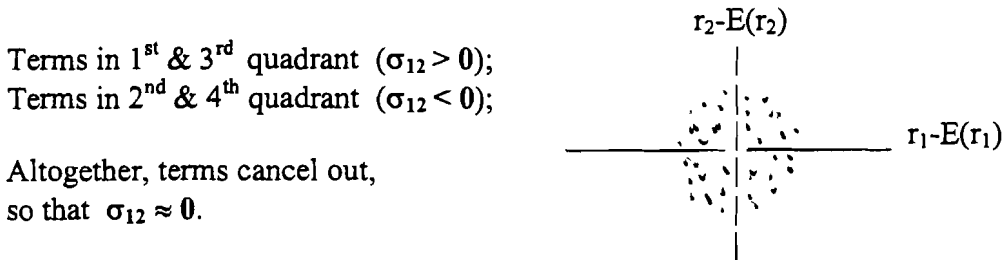
3. Case i: Suppose most combinations, (r_{1i}, r_{2i}) are in 1st & 3rd quadrants.



5. Case ii: Suppose most combinations, (r_{1i}, r_{2i}) are in 2nd & 4th quadrants.



6. Case iii: Suppose observations, (r_{1i}, r_{2i}) are scattered in all 4 quadrants.



Point: Sign of σ_{12} shows the *nature* of relation between r_1 & r_2 .

E. **Correlation** = $\rho_{12} = \sigma_{12} / \sigma_1 \sigma_2$.

1. Note: σ_{12} may vary between $-\infty$ & $+\infty$.
 - a. If r_1 and/or r_2 vary more widely (if σ_1 and/or σ_2 larger), then $[r_1 - E(r_1)]$ and/or $[r_2 - E(r_2)]$ are larger in magnitude, and σ_{12} will be larger in mag. (depending on case i, ii, or iii).
2. Thus, magnitude of σ_{12} doesn't tell us about extent of relation.
3. Correlation fixes this problem; adjusts σ_{12} for size of σ_1 and σ_2 .
 - a. If σ_{12} is larger (because σ_1 and/or σ_2 larger), ρ_{12} corrects for this by dividing σ_{12} by $(\sigma_1 \sigma_2)$.
4. Result: ρ_{12} varies between -1 and +1.
 - a. If $\rho_{12} = +1$, r_1 & r_2 are perfectly positively related.
 - b. If $\rho_{12} = -1$, r_1 & r_2 are perfectly negatively related.
 - c. If $\rho_{12} = 0$, r_1 & r_2 are unrelated.

Statistical Review [More Tools]

* [Throughout this explain with example.]

Let $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}$ be a parameter.

Defn 1: An estimator $\hat{\theta} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_p \end{bmatrix}$ of θ is any function of the observed sample.
(It depends on random variables.)

UNBIASED ESTIMATOR

Defn 2: $\hat{\theta}$ is unbiased if $E[\hat{\theta}] = \theta$.
i.e. if $E(\hat{\theta}_i) = \theta_i \quad \forall i=1, \dots, p$

Defn 3: The covariance matrix of $\hat{\theta}$ is:

$$\text{Cov}(\hat{\theta}) = E\{[\hat{\theta} - E(\hat{\theta})][\hat{\theta} - E(\hat{\theta})]'\}$$

note:

$$= E \begin{bmatrix} \hat{\theta}_1 - E(\hat{\theta}_1) \\ \hat{\theta}_2 - E(\hat{\theta}_2) \\ \vdots \\ \hat{\theta}_p - E(\hat{\theta}_p) \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 - E(\hat{\theta}_1) & \dots & \hat{\theta}_p - E(\hat{\theta}_p) \end{bmatrix}$$

→ p x p matrix { Variances on diag.
Covariances off dia

$$\text{Cov}(\hat{\theta})_{ij} = ij^{\text{th}} \text{ element}$$

$$= E\{[\hat{\theta}_i - E(\hat{\theta}_i)][\hat{\theta}_j - E(\hat{\theta}_j)]\}$$

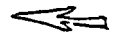
$$= \text{Cov}(\hat{\theta}_i, \hat{\theta}_j)$$

On diag.; $\text{Var}(\hat{\theta}_i)$

Off diag.; $\text{Cov}(\hat{\theta}_i, \hat{\theta}_j)$

EFFICIENT ESTIMATOR

Note: "Efficient" estimators are unbiased and have minimum variance (in some sense).



We want the $\hat{\theta}$ with min. covariance matrix

Defn 4: An unbiased estimator $\hat{\theta}$ is efficient relative to another unbiased estimator $\tilde{\theta}$ if

$$\text{Var}(\lambda' \hat{\theta}) \leq \text{Var}(\lambda' \tilde{\theta})$$

for every $p \times 1$ vector λ .

$$\text{Note: } \lambda' \hat{\theta} = [\lambda_1 \lambda_2 \dots \lambda_p] \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_p \end{bmatrix} = \lambda_1 \hat{\theta}_1 + \lambda_2 \hat{\theta}_2 + \dots + \lambda_p \hat{\theta}_p = \sum_{i=1}^p \lambda_i \hat{\theta}_i$$

Defn says that every l.c. of $\hat{\theta}$ has a variance at least as small as every l.c. of $\tilde{\theta}$.

→ Very strong defn! - \exists a whole lot of λ 's!

e.g. every element of $\hat{\theta}$ must have variance \leq every element of $\tilde{\theta}$.



Note: $\lambda' \hat{\theta}$ is a scalar; \therefore can talk about variance

Lemma: $\text{Var}(\lambda' \hat{\theta}) = \lambda' \text{Cov}(\hat{\theta}) \lambda$

yourself 3
 ← "Prove" with
 2x2 case

Proof: $\text{Var}(\lambda' \hat{\theta}) = \text{Var}\left[\sum_{i=1}^p \lambda_i \hat{\theta}_i\right]$

$= E\left\{\left[\sum_i \lambda_i \hat{\theta}_i - \sum_i \lambda_i E(\hat{\theta}_i)\right]^2\right\}$

$= E\left\{\sum_i \lambda_i [\hat{\theta}_i - E(\hat{\theta}_i)]\right\}^2$ pivot

$\lambda_i \neq \lambda_j$ not nec. =,
 just expanding
 the square (distributive property)

do yourself
 (circled)

$= E\left\{\sum_i \lambda_i [\hat{\theta}_i - E(\hat{\theta}_i)] \sum_j \lambda_j [\hat{\theta}_j - E(\hat{\theta}_j)]\right\}$

$= \sum_i \sum_j \lambda_i \lambda_j E[\hat{\theta}_i - E(\hat{\theta}_i)][\hat{\theta}_j - E(\hat{\theta}_j)]$

$= \sum_i \sum_j \lambda_i \lambda_j \text{Cov}(\hat{\theta}_i, \hat{\theta}_j)$

$= \lambda' \text{Cov}(\hat{\theta}) \lambda$ (quadratic form)

Thm:

An unbiased estimator $\hat{\theta}$ is efficient relative to an unbiased estimator $\tilde{\theta}$ iff

$[\text{Cov}(\tilde{\theta}) - \text{Cov}(\hat{\theta})]$ is p.s.d.

Proof: $\text{Var}(\lambda' \hat{\theta}) \leq \text{Var}(\lambda' \tilde{\theta}) \quad \forall \lambda$ ← (Defn 4) of efficiency

$\rightarrow \lambda' \text{Cov}(\hat{\theta}) \lambda \leq \lambda' \text{Cov}(\tilde{\theta}) \lambda \quad \forall \lambda$ by Lemma

$\rightarrow 0 \leq \lambda' [\text{Cov}(\tilde{\theta}) - \text{Cov}(\hat{\theta})] \lambda \quad \forall \lambda$

$\rightarrow [\text{Cov}(\tilde{\theta}) - \text{Cov}(\hat{\theta})]$ is psd (defn of psd)

Defn: An unbiased estimator $\hat{\theta}$ is efficient if it is efficient relative to all other unbiased estimators of θ .

* Thm (Cramer - Rao):

Let $\{x_1, x_2, \dots, x_n\}$ be a random sample (r.s.) from a population with density, $f(x; \theta)$.

Define the Likelihood function

$$L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

Define the Information Matrix, \mathcal{I} , by

$$\mathcal{I}_{ij} = - E \left[\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right] \quad i, j = 1, \dots, p$$

expected value

Then for any unbiased* estimator $\hat{\theta}$ of θ ,

$$[\text{Cov}(\hat{\theta}) - \mathcal{I}^{-1}] \text{ is psd.}$$

expl: \mathcal{I}^{-1} is the matrix of 2nd order partials of the log-likelihood function,

\mathcal{I}^{-1} is a "lower bound" on all estimates in $\text{Cov}(\hat{\theta})$ by this Thm.

① → correct: "lower bound" on all estimates on the diagonal of $\text{Cov}(\hat{\theta})$ (variance)

examples

Corollary: If $\text{Cov}(\hat{\theta}) = \mathcal{I}^{-1}$,
then $\hat{\theta}$ is efficient.

Note: This is a sufficient condition,
not necessary.

i.e. Can calculate $\text{Cov}(\hat{\theta})$ and \mathcal{I}^{-1} ,
if elements in $\text{Cov}(\hat{\theta}) =$ these in \mathcal{I}^{-1} ,
 $\hat{\theta}$ is efficient.

--- if \neq , $\hat{\theta}$ still may be efficient
(i.e. lower bound, \mathcal{I}^{-1} , may be unattainable)

Example:

$\{x_1, \dots, x_T\}$ is random sample from $N(\mu, \sigma^2)$

$$\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}; \quad f(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$$= (2\pi)^{-\frac{1}{2}} (\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$$\therefore L(x_1, \dots, x_T; \theta) = \prod_{i=1}^T (2\pi)^{-\frac{1}{2}} (\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\}$$

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

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

↳ Referred to a lot.

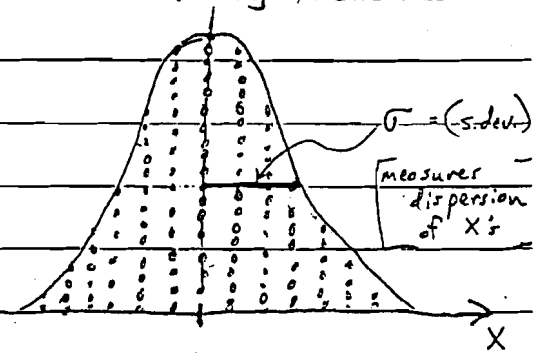
Intuitive expl. of Likelihood function.

eg. $\{x_1, \dots, x_T\}$ r.s. from $N(\mu, \sigma^2)$.

What does this mean?

- Draw out T x 's, observe their values. $f(x, \theta)$

If you were to plot a frequency distribution or histogram, these T obs. would pile up around μ . \rightarrow



Can draw a Bell-shaped curve around distribution.
 \rightarrow Normal distribution

Normal density is equation that gives this curve;

$$f(x; \theta) = (2\pi)^{-1/2} (\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2\sigma^2} (x-\mu)^2\right\}$$

Note: Only information we're given is x 's;
 don't know true μ and σ^2 .

- Want to estimate them from sample obs.!
 (Find out mean and variance of sample)

For each obs, x_i , can plug into $f(\theta; x_i)$
 and obtain numerical information about the distribution;
 presumably gives some point on the Bell-shaped curve.

--- \therefore get T values of density $f(\theta; x_i) \quad i=1 \quad T$

- These T values are like "probabilities" -
[the area under $f(x; \theta)$ is probability.]

How can we put all the information
contained in the sample together?

5. (INSERT B)

-- The Likelihood function puts all this
information together:

$$L = \prod_{i=1}^T f(\theta; x_i) = (2\pi)^{-\frac{T}{2}} (\sigma^2)^{-\frac{T}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^T (x_i - \mu)^2\right\}$$

This presents a number that depends on
the x_i 's and the parameter of interest, $\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}$;

$$L(x_1, \dots, x_T; \theta),$$

that gives us information about the distribution
the x_i 's are drawn from.

Presumably the $f(\theta; x_i)$, $i=1, \dots, T$,
somehow "trace out" the Bell-shaped curve,
and the Likelihood function incorporates all
this information [gives some "total probability" or "likelihood" about the dist]

Thus the Likelihood function provides information
that suggests to us what kind of distribution
it is most likely that these x_i 's would have
come from.

i.e. L tells us information about the distribution
(μ & σ^2) that would make it most likely that
we would observe what we have observed!

Note: form of L depends on kind of density the
sample is drawn from. Normal * Unimodal Dist. R. 1

24
P. 2
Handout

First Deriv.'s

$$\ln L = -\frac{1}{2} \ln(n!) - \frac{T}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

6

$$\frac{\partial \ln L}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_i 2(x_i - \mu)(-1) = \frac{1}{\sigma^2} \sum_i (x_i - \mu)$$

$$\frac{\partial \ln L}{\partial \sigma^2} = \frac{-T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_i (x_i - \mu)^2$$

Check if you want; $E[\]$ of each = 0.

(This regularity condition holds for

nearly all distributions.) [one exception; the uniform dist.]

Second Deriv.'s

$$\frac{\partial^2 \ln L}{\partial \mu^2} = \frac{1}{\sigma^2} \sum_i (-1) = \frac{-T}{\sigma^2}$$

$$\frac{\partial^2 \ln L}{\partial \sigma^2 \partial \mu} = \frac{\partial^2 \ln L}{\partial \mu \partial \sigma^2} = \frac{-1}{(\sigma^2)^2} \sum_i (x_i - \mu) = \frac{-1}{\sigma^4} \sum_i (x_i - \mu)$$

$$\frac{\partial^2 \ln L}{\partial (\sigma^2)^2} = \frac{T}{2(\sigma^2)^2} - \frac{2}{2(\sigma^2)^3} \sum_i (x_i - \mu)^2$$

$$= \frac{T}{2\sigma^4} - \frac{1}{\sigma^6} \sum_i (x_i - \mu)^2$$

Take Expected Values

$$E\left[\frac{\partial^2 \ln L}{\partial \mu^2}\right] = \frac{-T}{\sigma^2}$$

$$\begin{bmatrix} -T/\sigma^2 & 0 \\ 0 & -T/2\sigma^4 \end{bmatrix}$$

$$E\left[\frac{\partial^2 \ln L}{\partial \mu \partial \sigma^2}\right] = \frac{-1}{\sigma^4} E\left[\sum_i (x_i - \mu)\right] = 0 \quad (\text{since } E(x_i) = \mu)$$

$$E\left[\frac{\partial^2 \ln L}{\partial (\sigma^2)^2}\right] = \frac{T}{2\sigma^4} - \frac{1}{\sigma^6} E\left[\sum_i (x_i - \mu)^2\right] \quad (\text{since } E(x_i - \mu)^2 = \sigma^2)$$

$$= \frac{T}{2\sigma^4} - \frac{1}{\sigma^6} \sum_i \sigma^2 = \frac{T}{2\sigma^4} - \frac{T}{\sigma^6} \sigma^2 = \frac{-T}{2\sigma^4}$$

on board

Thus,

$$J = -E \begin{bmatrix} \frac{\partial^2 \ln L}{\partial \mu^2} & \frac{\partial^2 \ln L}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 \ln L}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 \ln L}{\partial (\sigma^2)^2} \end{bmatrix} = \begin{bmatrix} \frac{T}{\sigma^2} & 0 \\ 0 & \frac{T}{2\sigma^4} \end{bmatrix}$$

and

$$J^{-1} = \begin{bmatrix} \frac{\sigma^2}{T} & 0 \\ 0 & \frac{2\sigma^4}{T} \end{bmatrix}$$

← lower bound by Cramer's for $\text{Cov}(\hat{\theta})$

* Consider the estimator, $\hat{\theta} = \begin{bmatrix} \bar{X} \\ S^2 \end{bmatrix}$ of $\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}$.

where $\bar{X} = \frac{1}{T} \sum_i X_i$ and $S^2 = \frac{1}{T-1} \sum_i (X_i - \bar{X})^2$

and $\bar{X} \sim N(\mu, \frac{\sigma^2}{T})$ and $\frac{(T-1)S^2}{\sigma^2} \sim \chi_{T-1}^2$

Note: $E[\text{any } \chi^2] = \text{its d.f.}$

$\text{Var}[\text{any } \chi^2] = 2 * (\text{its d.f.})$

$$\text{Thus, } E\left[\frac{(T-1)S^2}{\sigma^2}\right] = T-1$$

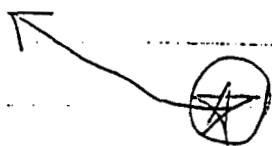
$$\rightarrow E\left[\frac{S^2}{\sigma^2}\right] = 1 \quad ; \quad \rightarrow E(S^2) = \sigma^2$$

$$\text{and } \text{Var}(S^2) = \text{Var}\left[\left(\frac{\sigma^2}{T-1}\right) \left(\frac{T-1}{\sigma^2}\right) S^2\right]$$

$$= \left(\frac{\sigma^2}{T-1}\right)^2 \text{Var}\left[\frac{(T-1)S^2}{\sigma^2}\right]$$

$$= \frac{\sigma^4}{(T-1)} = 2(T-1)$$

$$= \frac{2\sigma^4}{T-1}$$



Hence $E(\hat{\theta}) = E\left[\begin{matrix} \bar{X} \\ S^2 \end{matrix}\right] = \begin{bmatrix} \mu \\ \frac{\sigma^2}{T} \end{bmatrix} = \theta \rightarrow \underline{\text{unbiased}}$

and $\text{Cov}(\hat{\theta}) = \text{Cov}\left[\begin{matrix} \bar{X} \\ S^2 \end{matrix}\right] = \begin{bmatrix} \frac{\sigma^2}{T} & 0 \\ 0 & \frac{2\sigma^4}{T-1} \end{bmatrix}$ X_1 is nonstochastic
 $\therefore \text{Cov}(\bar{X}, S^2) = 0$

$[\text{Cov}(\hat{\theta}) - d^{-1}]$ is p.s.d. for any unbiased estimator!

Compare this with $d^{-1} : \begin{bmatrix} \sigma^2/T & 0 \\ 0 & 2\sigma^4/T \end{bmatrix}$

1) \bar{X} is efficient.

2) S^2 not proved efficient here;

Variance not as small as

lower bound.

$\text{Var}(\hat{\theta}_1)$

$\text{Var}(\hat{\theta}_2)$

efficiency measure

Consistency $\rightarrow \lim_{n \rightarrow \infty} P\{|\hat{\theta}_n - \theta| \leq \epsilon\} = 1$

$\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) = 0$

It turns out that

S^2 is efficient too.

We need more Thms to prove it.

[Note: The lower bound for S^2 in d^{-1} is unattainable here.]

OK next page

25

P. 3 Hands

Defn: Consider a random sample $\{x_1, \dots, x_T\}$ from $f(x; \theta)$. Let $t(x_1, \dots, x_T)$ be a function of the X 's.

Then t is sufficient for θ if the distribution of $\{x_1, \dots, x_T\}$ conditional on t does not depend on θ .

If $g_{x_1 \dots x_T | t}(x_1, \dots, x_T) \neq \theta$.

$L(x_1, \dots, x_T) = g(t, \theta) \cdot h(x_1, \dots, x_T)$

not a function of θ .

Intuition :

If t is sufficient, it contains all the information in the sample about θ , and it wouldn't make any difference if we were given (information about) θ or not.
(explicit)

* In other words, if t is a sufficient statistic for θ , it contains all the information w.r.t θ contained in the sample.

* Blackwell - Rao Thm :

A necessary and sufficient condition for $t(x_1, \dots, x_T)$ to be sufficient for θ is that

$$L(x_1, \dots, x_T; \theta) = g(t, \theta) \cdot h(x_1, \dots, x_T)$$

where g does not depend on $\{x_1, \dots, x_T\}$ except thru t , and h does not depend on θ .

→ i.e. You must be able to factor the likelihood function in this manner.

expl. with example! ↪

Example: $\{x_1, \dots, x_T\}$ is r.s. from $N(\mu, \sigma^2)$.

$$L = \left(\frac{1}{\sqrt{2\pi} \sigma}\right)^T \exp\left\{-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2\right\}$$

~~Try~~ Note: $\sum_i (x_i - \mu)^2 = \sum_i (x_i - \bar{x})^2 + T(\bar{x} - \mu)^2$

Substituting,

$$L = \left(\frac{1}{\sqrt{2\pi} \sigma}\right)^T \exp\left\{-\frac{1}{2\sigma^2} \left[\sum_i (x_i - \bar{x})^2 + T(\bar{x} - \mu)^2\right]\right\}$$
$$= \left(\frac{1}{\sqrt{2\pi} \sigma}\right)^T \exp\left\{-\frac{T}{2\sigma^2} \left[\bar{x} - \mu\right]^2 + \dots\right\}$$

Observe, $L =$ function of $\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}$ & $t = \begin{bmatrix} \bar{x} \\ S^2 \end{bmatrix}$

⊗ If $g =$ the entire function, L , and $h = 1$, then t is sufficient for θ .

Hence we have factored L into 2 functions;
 $g(t, \theta)$ which doesn't depend on $\{x_1, \dots, x_T\}$ except thru t ;

and $h(x_1, \dots, x_T)$ which doesn't depend on θ .



* The only place the x_i show up is in t , ($\bar{x} \neq \sigma^2$)

*** : (Min-variance unbiased estimator) Unbiased estimator that is a function of a minimal sufficient stat.

Note: Many other t 's work here;

e.g. $t' = \begin{bmatrix} \bar{x} \\ \sum (x_i - \bar{x})^2 \end{bmatrix}$, $t'' = \begin{bmatrix} \sum x_i \\ \sum x_i^2 \end{bmatrix}$, etc.

--- As long as the t 's contain all the information about θ in the sample, L can be factored into some $g(t, \theta) * \text{some } h$.

But also Note: If only one sufficient statistic, t , that is unbiased.

This brings up the following Thm.

Method of Lehmann & Scheffe: Say $(X_1 \dots X_n)$ and $(Y_1 \dots Y_n)$ comes from same dist.

If $\frac{L(X_1 \dots X_n)}{L(Y_1 \dots Y_n)}$ = free of θ then $g(Y_1 \dots Y_n)$ is the minimal sufficient stat for θ .

(on OH)

Thm: An unbiased estimator which is a function of a sufficient statistic is efficient.

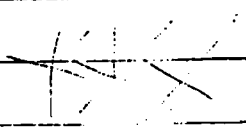
→ sufficient; uses all information

→ efficient; min. variance of all poss. estimators

Recall, showed that \bar{x} is efficient;

its variance = lower bound in σ^2

Now, know that s^2 is efficient by above Thm.



Sometimes must settle for less than an efficient estimator.

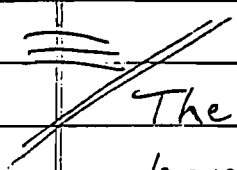
OH

Defn: An unbiased estimate of θ is the Best Linear Unbiased Estimator (BLUE) if

- ① it is a linear function of the sample
- and ② it is efficient relative to any other linear unbiased estimate.

Example: $\{x_1, \dots, x_n\}$ is r.s. from any distribution with finite mean and variance. Then \bar{x} is BLUE of μ .

Note: Already shown, if $\{x_1, \dots, x_n\}$ is Normal, \bar{x} is efficient estimator of μ .
 \Leftarrow (even beats nonlinear estimates)



The above statistical concepts have been dealing with "finite sample properties of estimators."

Now consider the Asymptotic (large sample) Properties of Estimators.

Asymptotic Properties hold only for ∞ samples,
But as $T \uparrow$, the estimators behave more closely
to their As. Properties.

13

Let $\hat{\theta}_T$ be an estimator of θ
based on a sample of T observations.

Useless Defn: $\hat{\theta}_T$ is asymptotically unbiased
if

$$\lim_{T \rightarrow \infty} E(\hat{\theta}_T) = \theta.$$

- Not very useful since we are often
interested in estimators where the moments
don't exist (e.g. Cauchy dist.).

ex. Consider $\hat{\theta}_T = \frac{T-1}{T} \bar{X}$. - Biased (est of μ),
but As. Unbiased.

More
Useful

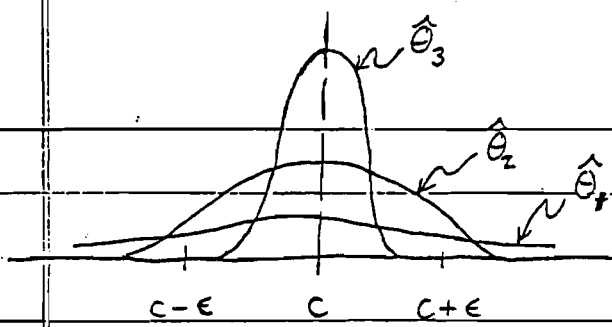
Defn: Consider a sequence of random variables
e.g. $\{\hat{\theta}_1, \hat{\theta}_2, \dots\}$
Then if \exists a number c \int

$$\lim_{T \rightarrow \infty} P(|\hat{\theta}_T - c| > \epsilon) = 0 \quad \forall \epsilon > 0$$

then the sequence converges in probability
to c , and we write

$$p\lim_{T \rightarrow \infty} \hat{\theta}_T = c.$$

expl. \int



The probability (area under pdf's, outside of $c \pm \epsilon$) approaches 0 as $T \rightarrow \infty$.

..... The distribution of the $\hat{\theta}$'s converges to a "spike" at c .
 -- The variance of $\hat{\theta}_T \rightarrow 0$ as $T \rightarrow \infty$. (★)

(Note:) the above 2 Defn's ^{almost} say the same thing. The second is more useful since we don't need to take any expected values to make use of it.
 - plim also \rightarrow variance goes to zero!

Example: $\{z_1, z_2, \dots, z_T\}$ a r.s. from $N(\mu, \sigma^2)$

$$\bar{z}_T = \frac{1}{T} \sum z_i \quad ; \quad \bar{z}_T \sim N\left(\mu, \frac{\sigma^2}{T}\right)$$

The variance $\rightarrow 0$ as $T \rightarrow \infty$ } ↗

$$\lim_{T \rightarrow \infty} P(|\bar{z}_T - \mu| > \epsilon) = 0,$$

$$\text{or } \text{plim } \bar{z}_T = \mu.$$

This ex. reflects the Law of Large Numbers, and shows the As. Properties of \bar{z}_T .

Defn: An estimator $\hat{\theta}_T$ is consistent if $\text{plim}_{T \rightarrow \infty} \hat{\theta}_T = \theta$.

— implies 2 conditions:

- 1) $E(\hat{\theta}_T) \rightarrow \theta$
- 2) $\text{Var}(\hat{\theta}_T) \rightarrow 0$

Defn: The Mean Square Error of $\hat{\theta}_T$ is

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

and note ...

$$\begin{aligned}
 &= E[(\hat{\theta}_T - E(\hat{\theta}_T)) + (E(\hat{\theta}_T) - \theta)]^2 \\
 &= E[\hat{\theta}_T - E(\hat{\theta}_T)]^2 + [E(\hat{\theta}_T) - \theta]^2 + \\
 &\quad + 2[E(\hat{\theta}_T) - \theta] \underbrace{E[\hat{\theta}_T - E(\hat{\theta}_T)]}_0 \\
 &= E[\hat{\theta}_T - E(\hat{\theta}_T)]^2 + [E(\hat{\theta}_T) - \theta]^2 \\
 &= \underline{\underline{\text{Var}(\hat{\theta}_T)}} + \underline{\underline{[\text{Bias}]^2}}
 \end{aligned}$$

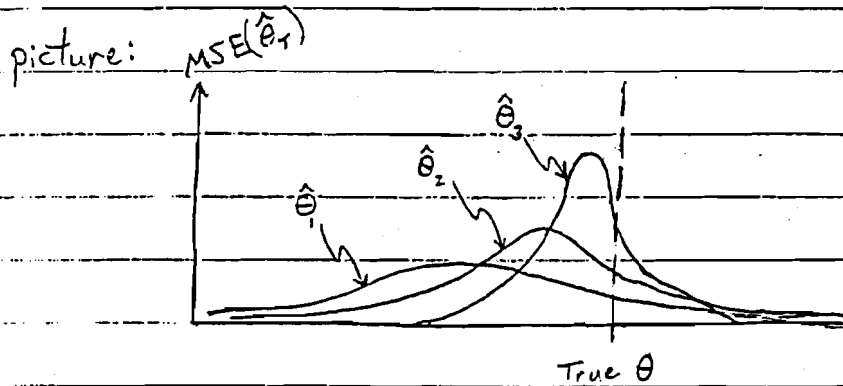
Note: The $\text{MSE}(\hat{\theta}_T) = \text{Var}(\hat{\theta}_T)$ if $\hat{\theta}_T$ is unbiased.

Thus $\text{MSE}(\hat{\theta}_T)$ is a criterion used to judge how good an estimator is, if it is biased.
— just like $\text{Var}(\hat{\theta}_T)$ is used to judge unbiased estimates.

Thm: IF $\lim_{T \rightarrow \infty} \text{MSE}(\hat{\theta}) = 0$,

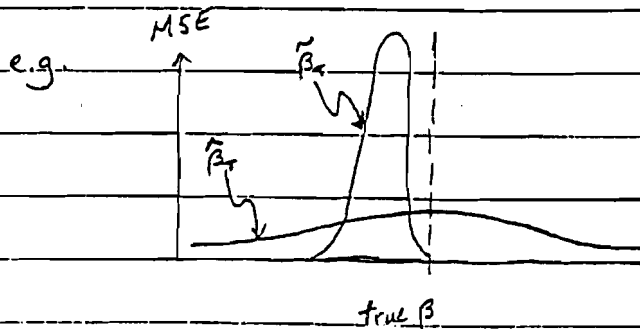
then $\hat{\theta}_T$ is consistent.

Proof: meets 2 conditions -
Variance & Bias go to zero.



Note: "Bias" is not always a terrible monster to avoid at all costs.

In some cases a Biased estimate $\tilde{\beta}$ may be preferable to an unbiased est. $\hat{\beta}$



→ here $\text{MSE}(\tilde{\beta}) < \text{MSE}(\hat{\beta})$

[note here that $\tilde{\beta}$ will never yield true β !!]

— Better to get close to true β with high degree of precision, than to have unbiased estimate with huge variance. (low precision)

Thm: (Slutsky)

If $\text{plim } x_T = c$, and if g is any continuous function, then

$$\text{plim } g(x_T) = g(\text{plim } x_T) = g(c)$$

Examples: ① $\text{plim } (\hat{\theta}_T)^2 = (\text{plim } \hat{\theta}_T)^2$

$$\text{② } \hat{\theta}, \tilde{\theta}; \text{plim} \left(\frac{\hat{\theta}}{\tilde{\theta}} \right) = \frac{\text{plim } \hat{\theta}}{\text{plim } \tilde{\theta}}$$

Note: This does not hold in general for Expected Values.

$$\text{e.g. } [E(x)]^2 \neq E(x^2)$$

→ unless $\text{Var}(x) = 0$

↳ i.e. $E(x^2) = [E(x)]^2 + \text{Var}$

Discussion:

① The plim is the (large sample) counterpart to expected value wrt finite samples.

② Many textbook treatments of asymptotic distributions are wrong because they frequently use expectations (moments), which do not always exist.

verify \longrightarrow (e.g. Johnston, Kmenta)

Defn: Let $\hat{\theta}_T$ be a consistent estimator of θ \downarrow
 $\sqrt{T}(\hat{\theta}_T - \theta)$ converges in distribution to $N(0, \Omega)$.
 Then the Asymptotic Distribution of $\hat{\theta}$ is $N(\theta, \frac{1}{T}\Omega)$.

(where Ω is the covariance matrix of $\hat{\theta}$.)

Intuitive expl.:

** $(\hat{\theta}_T - \theta) \rightarrow 0$ as fast as \sqrt{T} changes as $T \rightarrow \infty$;

--- Instead of the probability density functions converging to a point (as with plim's),

their cumulative distribution functions converge to a known or fixed cum. dist. function

Shady Comment:

Given $E(\hat{\theta}_T) \rightarrow \theta$ and $\text{Cov}(\hat{\theta}_T) \rightarrow \frac{1}{T}\Omega$,
 it is often true that

$$\frac{\hat{\theta}_T - \theta}{\frac{\sqrt{\Omega}}{\sqrt{T}}} \rightarrow N(0, 1)$$

— by Central Limit Thm.

Explain with examples!

Trivial

Example: $\{x_1, \dots, x_T\}$ is a r.v. from any distribution, with mean μ and variance σ^2 .

Consider $\bar{x}_T = \frac{1}{T} \sum_i x_i$; $E(\bar{x}) = \mu$, $Var(\bar{x}) = \frac{\sigma^2}{T}$.

$\text{plim } \bar{x} = \mu$ by the Law of Large Numbers (shown before) (\bar{x} consistent)

Then $\frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{T}}} \rightarrow N(0, 1)$ by Central Limit Thm.

Note: $\frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{T}}} = \frac{\sqrt{T}(\bar{x} - \mu)}{\sigma} \rightarrow N(0, 1)$.

☆☆ Thus $\sqrt{T}(\bar{x} - \mu) \rightarrow N(0, \sigma^2)$.

Then it follows from our defn, that

$$\bar{x} \rightarrow N\left(\mu, \frac{1}{T} \sigma^2\right),$$

⊥ we know this is true even for small samples.

" \sqrt{T} " is like a magic normalization factor.

Normally, wrt this example, we don't talk about \bar{x} converging in distribution to this because

this dist. is a $f(\frac{1}{T})$ ⊥ this is messy as $T \rightarrow \infty$. Typically ~~normalize~~ ^{standardize} \bar{x} ⊥ then work with $N(0, 1)$ which is not a $f(1/T)$

Useful Fact:

~~XXXXXX~~ A statistic $\sim N(\mu, \sigma^2)$ multiplied by a constant, c , is distributed $N(c\mu, c^2\sigma^2)$

i.e. $c * N(\mu, \sigma^2) \sim N(c\mu, c^2\sigma^2)$

Now.

A less trivial example.

$\{X_1, \dots, X_T\}$ a r.s. from $N(\mu, \sigma^2)$

$S^2 = \frac{1}{T-1} \sum_i (X_i - \bar{X})^2$ is a consistent estimator of σ^2 .

Recall $\frac{(T-1)S^2}{\sigma^2} \sim \chi^2_{T-1}$ Mean: $T-1$ Var: $2(T-1)$

Consider

(Algebra) \rightarrow
(Take $\lim_{T \rightarrow \infty}$) \rightarrow

$$\sqrt{T}(S^2 - \sigma^2) = \underbrace{\sqrt{2}\sigma^2}_{\sqrt{2}\sigma^2} \underbrace{\sqrt{\frac{T}{T-1}}}_{(1)} \underbrace{\left[\frac{\frac{(T-1)S^2}{\sigma^2} - (T-1)}{\sqrt{2(T-1)}} \right]}_{N(0,1)}$$

Hence, $\sqrt{T}(S^2 - \sigma^2) \rightarrow \sqrt{2}\sigma^2 * N(0, 1)$

or $\sqrt{T}(S^2 - \sigma^2) \rightarrow N(0, 2\sigma^4)$ [by useful fact]

or The Asymptotic Dist. of S^2 is $N(\sigma^2, \frac{2\sigma^4}{T})$

[by Defn, p. 18]

* Nice to know!

As. Variance of S^2 is Cramer-Rao lower bound (1-7) \rightarrow

Useful Fact Restated: (more general form)

If A is a matrix $\} \text{plim } A \text{ exists,}$
 $\& \text{ if } \xi \text{ is a vector } \} \xi \rightarrow N(0, \Omega),$

then $A\xi \rightarrow N[0, (\text{plim } A)\Omega(\text{plim } A)']$.

Note: $\text{plim } A \Rightarrow \text{plim of each element in}$

In this example,

let $A = \sqrt{2} \sigma^2 \sqrt{\frac{T}{T-1}}$; $\xi = \text{our standardized } \chi^2 \text{ statistic;}$

then from the Algebra, $\sqrt{A\xi} = \sqrt{T}(\sigma^2 - \sigma^2)$, ^{previous} page !!!

and $A\xi \rightarrow N[0, (\text{plim } A)\Omega(\text{plim } A)']$

or $\rightarrow N[0, (\sqrt{2} \sigma^2) 1 (\sqrt{2} \sigma^2)]$

or $\rightarrow N(0, 2\sigma^4)$.

~~* Also, $\text{plim } A = \sqrt{2} \sigma^2$~~

Defn:

If $\hat{\theta} \neq \tilde{\theta}$ are consistent estimators of θ with asymptotic covariance matrices $\frac{1}{T}\Omega \neq \frac{1}{T}\Psi$, respectively, then

$\hat{\theta}$ is asymptotically efficient relative to $\tilde{\theta}$ if $(\Psi - \Omega)$ is psd.

----(analogous to finite sample defn)

Defn: A consistent estimator $\hat{\theta}$ is asymptotically efficient if it is asymptotically efficient relative to every other consistent estimator of θ .

Defn: Let $\{x_1, \dots, x_T\}$ be a r.s. from a population with density, $f(x; \theta)$. Then the Maximum Likelihood Estimator (MLE) of θ is the value that maximizes the Likelihood function, $L = \prod_{t=1}^T f(x_t; \theta)$, wrt θ .

Intuition: Picking the value of $\hat{\theta}$ that makes the sample observations most likely;
i.e. the value of $\hat{\theta}$ that makes it most likely that you would observe what you have observed!

Example: $\{x_1, \dots, x_T\}$ is a r.s. from $N(\mu, \sigma^2)$.

$$\text{Recall, } L = (2\pi)^{-T/2} (\sigma^2)^{-T/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2\right\}$$

$$\ln L = -\frac{T}{2} \ln(2\pi) - \frac{T}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2$$

----- maximize this wrt μ & σ^2 .

Note: since the \ln function is monotonic,
the values that max. $(\ln L)$ also max. (L) .

$$\frac{d \ln L}{d \mu} = \frac{1}{\sigma^2} \sum_i (x_i - \hat{\mu}) = 0$$

$$\frac{d \ln L}{d \sigma^2} = \frac{-T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_i (x_i - \hat{\mu})^2 = 0$$

Convention: When we set 1st deriv.'s = 0,
the parameters become MLE's,
so put little hats on them ($\hat{\cdot}$).

$$\text{Solving, we get: } \hat{\mu} = \bar{x} \\ \hat{\sigma}^2 = \frac{1}{T} \sum_i (x_i - \bar{x})^2 \quad \left[= \left(\frac{T-1}{T}\right) s^2 \right]$$

Very Important

Thm: Under general conditions*,
the MLE $\hat{\theta}$ of a parameter θ
is consistent and asymptotically efficient.
Furthermore,

$$\sqrt{T} (\hat{\theta} - \theta) \rightarrow N\left[0, \lim_{T \rightarrow \infty} \left(\frac{dL}{dT}\right)^{-1}\right]$$

*See Appendix of Schmidt.

This Thm means that
the Asymptotic Covariance Matrix of $\hat{\theta}$ is

$$\frac{1}{T} \lim_{T \rightarrow \infty} \left(\frac{d}{T} \right)^{-1} = \frac{1}{T} \lim_{T \rightarrow \infty} T(d^{-1})$$

$$\downarrow$$

$$= d^{-1} \quad \left[\text{if } T \text{ is sufficiently large.} \right]$$

-- Simply the Information Matrix!
Hence MLE's are asymptotically efficient.

Thus we'll often work with MLE's --- they're good!

--- Many computer programs present information about the Likelihood function, e.g. $\ln L$.

MLE's will be the estimates that maximize this.

Hence, can change model as desired,

& choose the one that reaches a higher point on the Likelihood surface...

Thm: Let $\hat{\theta}$ be any consistent estimator of θ .
Define the Linearized MLE as

(Convert inefficient estimator, $\hat{\theta}$ to efficient one $\tilde{\theta}$)

$$\tilde{\theta} = \hat{\theta} - \underbrace{\left[\frac{\partial^2 \ln L}{\partial \theta \partial \theta'} \right]^{-1}}_{\substack{\text{p} \times \text{p} \\ \text{(similar to } d)}} \left[\frac{\partial \ln L}{\partial \theta} \right]$$

$\begin{matrix} \text{p} \times 1 & \text{p} \times 1 & & \text{p} \times 1 \end{matrix}$

with all derivatives evaluated at $\hat{\theta}$.

Then $\tilde{\theta}$ is consistent and asymptotically efficient

This is a way of converting an asymptotically inefficient estimator into an asymptotically efficient one. ($\hat{\theta}$ into $\tilde{\theta}$)

Called the Newton-Raphson Algorithm

or the Method of Scoring.

Used in computer programs;
saves the more costly work of
evaluating all the MLE's.

[Intuition: get functional form of
1st & 2nd derivatives in the Algorithm;
plug in values of $\hat{\theta}$ to get $\tilde{\theta}$.

It is good if $\hat{\theta}$ is close to begin with,
not good if not...

$\hat{\theta}_T$ is asymptotically unbiased
if

$$\lim_{T \rightarrow \infty} E(\hat{\theta}_T) = \theta.$$

25. (ISERT.A)

Defn: Consider a sequence of random variables
eg. $\{\hat{\theta}_1, \hat{\theta}_2, \dots\}$.

Then if \exists a number c s.t.

$$\lim_{T \rightarrow \infty} P(|\hat{\theta}_T - c| > \epsilon) = 0 \quad \forall \epsilon > 0$$

then the sequence converges in probability
to c , and we write

$$\text{plim } \hat{\theta}_T = c.$$

Defn: An estimator $\hat{\theta}_T$ is consistent if
 $\text{plim } \hat{\theta}_T = \theta$.

Defn: The Mean Square Error of $\hat{\theta}_T$ is

$$\text{MSE}(\hat{\theta}_T) = E(\hat{\theta}_T - \theta)^2.$$

Thm: If $\lim_{T \rightarrow \infty} \text{MSE}(\hat{\theta}_T) = 0$,

then $\hat{\theta}_T$ is consistent.

Thm: (Slutsky)

If $\text{plim } X_T = c$, and if g is any
continuous function, then

$$\text{plim } g(X_T) = g(\text{plim } X_T) = g(c)$$

Consider 2 Well-known Tests of Hypotheses.

i) Likelihood Ratio Test (Neyman-Pearson Application)

← Consider some function of θ , $f(\theta) = 0$. $\leftarrow H_0$ (null hypothesis) $(m \times 1)$

a) Maximize L wrt θ ;
observe $L_u =$ maximized value of L , unrestricted.

$$L = L(X_1, \dots, X_T; \theta = \hat{\theta} = \text{MLE of } \theta) \quad \text{unrestr.}$$

b) Maximize L ~~unrestricted~~ s.t. $f(\theta) = 0$;
observe $L_r =$ maximized value of L , restricted.
--- [get MLE of θ s.t. $H_0: \theta = c$]

Obviously $L_u \geq L_r$.

The magnitude of their difference indicates whether the hypothesis that $f(\theta) = 0$ [H_0] can be accepted or rejected.

If $H_0: f(\theta) = 0$ is true,

L_u should approximately = L_r ;

i.e. the constraint should not be severely binding.

↳ Under H_0 , $-2 \ln\left(\frac{L_r}{L_u}\right) \rightarrow \chi^2_{lm}$

Reject H_0 if $-2 \ln\left(\frac{L_r}{L_u}\right) = 2(\ln L_u - \ln L_r)$ is large.

ii) The Wald Test

Let $\hat{\theta}$ = the unrestricted MLE of θ .

As a test statistic, consider $f(\hat{\theta})$;
is it close to 0?

$H_0: \theta = \alpha$
 $f(\theta) = \theta - \alpha$

Suppose we know the asymptotic dist. of $\hat{\theta}$:

$$\sqrt{T} f(\hat{\theta}) \rightarrow N(0, \Omega)$$

It's not
"Suppose"

It comes from
the property of MLE.

$$\sqrt{T} (\hat{\theta} - \theta) \rightarrow N(0, \frac{\text{Var}(\frac{\partial \log L}{\partial \theta})}{T})$$

[i.e. suppose we can prove the
Central Limit Theorem in this case.]

Then

$$T [f(\hat{\theta})' \Omega^{-1} f(\hat{\theta})] \rightarrow \chi^2_m$$

$1 \times m \quad m \times m \quad m \times 1$

Under H_0

$f(\hat{\theta})$

* Same dist. as Likelihood Ratio Test.

* In calculating this test statistic,
evaluate Ω^{-1} at $\hat{\theta}$ (where-ever θ appears).

* Advantage: doesn't require construction of ~~restricted~~
constrained MLE.

* See Judge, p.757

EXAM 2

Consider a random sample $\{X_1, \dots, X_T\}$

from a population with density

$$f(X) = \frac{1}{\theta} e^{-\frac{1}{\theta} X} \quad \text{for } X \geq 0$$

- 3
- A. Find the Maximum Likelihood Estimator of θ .
- 1
- B. Prove that it is unbiased.
- 5
- C. Prove that it is efficient.
- 1
- D. How would you test the null hypothesis,

$$H_0: \theta - 3 = 0 \quad ?$$

Useful Fact: The mean of the above distribution is θ ,
and the variance is θ^2 .