

INTRODUCTION AND BACKGROUND RESULTS

"Moments" of a Distribution:

12.

$y \in Y$
 $\theta \in \Omega$
 (sample space)

$p(y|\theta)$
 (parameter space)
 (data density)

$S_n = S(y_1, \dots, y_n)$
 (statistic)

The probability distribution of S_n is called
 its "sampling distribution".

Our task is to use the data, via S_n , to
 draw inferences about θ .

Estimators and tests are decision rules.

The characteristics of the sampling dist'n.
 of S_n will essentially determine its
 quality as a basis for inference.

e.g.: $y_i \sim (\mu, \sigma^2)$ (i.i.d.)

$S_n = \bar{y} = \frac{1}{n} \sum y_i \sim (\mu, \sigma^2/n)$

Let Y be a random variable with
 density $p(y|\theta)$. Then the "moments" of
 Y are defined as:

$$(i) E(Y^k) = \mu'_k = \int_{-\infty}^{\infty} y^k p(y|\theta) dy ; k=1, 2, \dots$$

("Raw moments" or "moments about origin")

$$(ii) E(Y - \mu'_1)^k = \mu_k = \int_{-\infty}^{\infty} (y - \mu'_1)^k p(y|\theta) dy \quad (k=1, 2, \dots)$$

("moments about the mean")

$$\text{So: } \mu'_1 = E(Y) = \text{mean}$$

$$\begin{aligned} \mu'_2 &= E(Y - \mu'_1)^2 = \text{variance} \\ &= E(Y^2) + (\mu'_1)^2 - 2\mu'_1 E(Y) \\ &= E(Y^2) - (\mu'_1)^2 \\ &= \mu''_2 - (\mu'_1)^2 \end{aligned}$$

We can always write the moments
 about mean in terms of moments about
 zero, & vice versa.

Why are moments important?

If all of the moments exist, then knowledge of the moments provides knowledge of the underlying distribution.

If all of moments exist, two distributions

will be identical if their moments match.

Uniqueness requires existence of moments).

There is a convenient way of generating the moments of a distribution, without

any need to integrate. We take the

Fourier transform of the density, which is called the Characteristic Function of the random variable.

Definition: The characteristic function of γ

$$\phi_{\gamma}(t) = E(e^{it\gamma}) ; t^2 = -1$$

$$= \int_{-\infty}^{\infty} e^{ity} p(y) dy$$

Example:

$$\gamma \sim N(0, 1)$$

$$\begin{aligned} \text{So, } \phi_{\gamma}(t) &= E(e^{it\gamma}) = \int_{-\infty}^{\infty} e^{ity} \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(y+it)^2/2} \cdot e^{-t^2/2} dy \\ &= e^{-t^2/2} \quad (= e^{-(it)^2/2}) \end{aligned}$$

Show that:

$$(i) \quad \text{If } \gamma \sim N(\mu, \sigma^2), \text{ then } \phi_{\gamma}(t) = \exp[i\mu t - \frac{\sigma^2 t^2}{2}] .$$

(ii) If $\gamma_1 + \gamma_2$ are independent, then the c.f. of their sum is the product of their individual c.f.'s.
 $[\phi_{\gamma_1+\gamma_2}(t) = \phi_{\gamma_1}(t) \cdot \phi_{\gamma_2}(t)]$

Example: Let $z \sim N(0, 1)$ & $\gamma = z^2$.

$$\phi_{\gamma}(t) = \int_{-\infty}^{\infty} e^{itz^2} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz$$

$$\begin{aligned} &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-z^2/2(1-2it)} dz \end{aligned}$$

$$= (1-2it)^{-k} \int_{-\infty}^{\infty} \frac{1}{(1-2it)^{-k}\sqrt{2\pi}} e^{-\frac{z^2}{2(1-2it)^{-1}}} dz$$

$$= (1-2it)^{-1/2}$$

So, the c.f. for $X_{(p)}^2$ is $(1-2it)^{-1/2}$.

Similarly, the c.f. for $X_{(p)}^4$ is $(1-2it)^{-p/2}$.

There are various ways of getting the moments of the distribution from the c.f.:

(a) Differentiation -

$$\mu'_k = \left[\frac{d^k \phi(t)}{dt^k} \right] / t^{k+0}$$

e.g. : $\gamma \sim \chi^2_{(p)}$; $\phi_\gamma(t) = (1-2it)^{-p/2}$

$$\phi'_\gamma(t) = (-pt)(-2i)(1-2it)^{-p/2-1}$$

$$= i\gamma (1-2it)^{-(p/2+1)}$$

$$\therefore \mu'_1 = [\phi'_\gamma(t)|_{t=0}] / i = p.$$

$$\phi''_\gamma(t) = -(\gamma p)(pt+1)(-2i)(1-2it)^{-(p/2+2)}$$

$$\therefore \mu'_2 = [\phi''_\gamma(t)|_{t=0}] / i^2 = 2p(p/2+1)$$

$$= 2p(p+2)/2 = p(p+2)$$

$$\text{So, var.}(\gamma) = \mu_2 = \mu'_2 - (\mu'_1)^2 = 2p.$$

(b) Taylor-Series Expansion -

μ'_k = coefficient of $(it)^k/k!$ in series expansion.

e.g. $\gamma \sim N(0, 1)$; $\phi_\gamma(t) = e^{-t^2/2}$

$$\text{Now, } e^x = 1 + x + \frac{x^2}{2!} + \dots$$

$$\text{So, } \phi_\gamma(t) = 1 + \frac{(it)^2}{2} + \left[\frac{(it)^2}{2} \right]^2 / 2! + \left[\frac{(it)^2}{2} \right]^3 / 3! + \left[\frac{(it)^2}{2} \right]^4 / 4! + \dots$$

$$= 1 + \frac{(it)^2}{2!} + \frac{(it)^4}{8} + \frac{(it)^6}{48} + \dots$$

So, $\mu'_k = 0$; for all odd k .

$$\mu'_2 = 1$$

$$\mu'_3 = 3 \quad \left(\frac{1}{8} = 3/4! \right)$$

$$\text{i.e. : E}(\gamma) = 0 \quad ; \quad \text{var.}(\gamma) = 1 - 0^2 = 1.$$

Definition : Skewness = $\frac{\mu_3}{(\mu_2)^{3/2}}$

$$\therefore \mu_3 = E(\gamma - \mu)^3 = \mu'_3 + 2(\mu'_1)^3 - 3\mu'_1\mu'_2$$

e.g. $Y \sim N(0, 1)$

$$\mu_1' = \mu_3' = 0 ; \mu_2' = 1 ; \mu_4' = 2$$

So, Skewness = 0

Definition: Kurtosis = $\left[\frac{\mu_4}{\mu_2^2} \right]$

$$\begin{aligned} \therefore \mu_4 &= E(Y - \mu)^4 = \mu_4' - 4\mu_3'\mu_1' + 6\mu_2'\mu_1'^2 \\ &\quad - 3(\mu_1')^4 \end{aligned}$$

e.g. $Y \sim N(0, 1)$

$$\mu_4 = 3 - 0 + 0 - 0 = 3$$

(Also true for $N(\mu, \sigma)$)

Normal Distn. is "Mesokurtic."

If kurtosis > 3 ,

"Leptokurtic"

e.g. Student 't'

If kurtosis < 3 ,

"Platykurtic"

e.g. Uniform

So, we often report the "excess kurtosis":
$$\left[\frac{\mu_4}{\mu_2^2} \right] - 3$$

Summary :

- * Moments of a distribution fully describe that distn. (if they exist).
- * Characteristic func. provides a painless way of obtaining the moments.

The moments are used to determine

properties of a statistic when used as estimator or test statistic.

Estimator Properties :

Focus on "exact" (finite-sample) properties. Consider asymptotic properties after introducing Maximum Likelihood estimation.

(i) Bias:
 $Bias(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta$

("Unbiased" if $Bias = 0$)
Compares first moment of $\hat{\theta}_n$ with location

(ii) (Relative) Efficiency:

Let $\hat{\theta}_n$ & $\tilde{\theta}_n$ be 2 unbiased estimators of θ . Then $\hat{\theta}_n$ is efficient relative to $\tilde{\theta}_n$ if $[V(\tilde{\theta}_n) - V(\hat{\theta}_n)]$ is p.s.d.

(Scalar : var. ($\tilde{\theta}_n$) \geq var. ($\hat{\theta}_n$))

If one or both of the estimators is biased, $\hat{\theta}_n$ is relatively more efficient than $\tilde{\theta}_n$ if $[M(\tilde{\theta}_n) - M(\hat{\theta}_n)]$ is p.s.d.,

where $M(\theta_n^*) = V(\theta_n^*) + \text{Bias}(\theta_n^*)\text{Bias}(\theta_n^*)'$
is the matrix MSE. (1st. & 2nd. moments)

(iii) Sufficiency:

If we can find a statistic, \bar{s}_n , that tells us as much about the population parameter(s) as does the full sample, this statistic will be "sufficient" for estimation purposes. More formally —

9

40

Let $\{y_i\}_{i=1}^n$ be a random sample from $P(y|\theta)$. Then \bar{s}_n is a "sufficient statistic" iff $P(y|s_n)$ does not depend on θ .

(If you know s_n , the sample values themselves add no information about θ .)

This idea extends to "jointly sufficient statistics". And, any 1-1 transformation of a set of jointly sufficient statistics is also jointly sufficient.

Example:

If \bar{z}_n & \bar{y}_n are jointly sufficient, then so are \bar{y} and $\bar{z}_n(y_i - \bar{y})^2$.

Now, how do we find a sufficient statistic in practice?

Factorization Theorem:

$\sum y_i$ is sufficient iff

$$p(y_1, \dots, y_n | \theta) = g(\sum y_i | \theta) h(y_1, \dots, y_n)$$

where $h(\cdot) \geq 0$ & does not involve θ ,
and $g(\cdot) \geq 0$ & depends on the y_i 's only
through $\sum y_i$.

For jointly sufficient statistics, this
becomes:

$$p(y_1, \dots, y_n | \theta) = g(\sum y_i, \dots, \sum y_n | \theta) h(y_1, \dots, y_n)$$

Example: Let y_1, \dots, y_n be drawn randomly
from a Bernoulli density:

$$\begin{aligned} p(y_1 | \theta) &= \prod_{i=1}^n p(y_i | \theta) \\ &= \prod_{i=1}^n \theta^{y_i} (1-\theta)^{1-y_i}; \quad 0 \leq \theta \leq 1 \\ &= \theta^{\sum y_i} (1-\theta)^{n-\sum y_i} \end{aligned}$$

$$\text{So, } g(\cdot) = e^{\sum y_i} (1-\theta)^{n-\sum y_i}; \quad h(\cdot) = 1$$

& $\sum y_i$ is sufficient for θ (so is \bar{y})

Example:

$$\begin{aligned} y_i &\sim \text{iid } N(\mu, \sigma^2) \\ f(y | \mu) &= \frac{1}{\sqrt{2\pi}} \exp[-\frac{1}{2}(y_i - \mu)^2] \\ &= \frac{1}{(2\pi)^{n/2}} \exp[-\frac{1}{2} \sum (y_i - \mu)^2] \\ &= \frac{1}{(2\pi)^{n/2}} \exp[-\frac{1}{2} (\sum y_i^2 + 2\mu \sum y_i + n\mu^2)] \\ &= \frac{1}{(2\pi)^{n/2}} \exp[\mu \sum y_i - \frac{1}{2}\mu^2] \exp[-\frac{1}{2} \sum y_i^2] \end{aligned}$$

$$\text{Let } g(\cdot) = \exp[\mu \sum y_i - \frac{1}{2}\mu^2]$$

$$h(\cdot) = \frac{1}{(2\pi)^{n/2}} \exp[-\frac{1}{2} \sum y_i^2]$$

Then, $\sum y_i$ is a sufficient statistic,
& so is \bar{y} .

Intuitively, it looks as if sufficient stats.
(or simple funcs. of them) may lead to
familiar & sensible estimators - this
will indeed turn out to be the case.

Now let's pursue the idea of estimation as a decision-making procedure & introduce some ideas from statistical decision theory.

(iv) Risk:

Suppose we have a loss function,

$$L(\hat{\theta}_n; \theta) \geq 0 ; \text{ all } \theta, \hat{\theta}_n$$

$$= 0 ; \text{ if } \hat{\theta}_n = \theta.$$

(N.B. : Loss is random)

To get a single value, average :

$$\text{Risk}(\hat{\theta}_n; \theta) = R(\hat{\theta}_n; \theta)$$

$$= E[L(\hat{\theta}_n; \theta)]$$

$$= \int L(\hat{\theta}_n; \theta) P(y|\theta) dy$$

e.g. Quadratic Loss -

$$(i) \text{ Scalar : } L = c(\hat{\theta}_n - \theta)^2 ; c > 0$$

$R = \text{MSE}$, if $c=1$.

$$(ii) \text{ Vector : } L = (\hat{\theta}_n - \theta)' W (\hat{\theta}_n - \theta) ; W \text{ pd.}$$

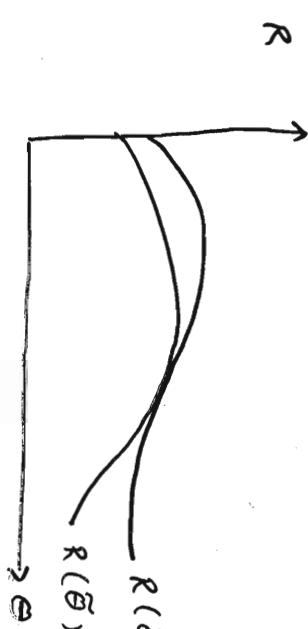
$$R = \text{tr}(W \text{MSE}) , \text{ if } W = I.$$

13.

(v) Admissibility -

Low risk seems desirable. $\hat{\theta}_n$, an estimator (decision rule), is Inadmissible if \exists any $\tilde{\theta}_n$ such that $R(\tilde{\theta}_n) \leq R(\hat{\theta}_n)$ for all θ , and $R(\tilde{\theta}_n) < R(\hat{\theta}_n)$ for some θ .

An Admissible estimator is one which is not inadmissible.



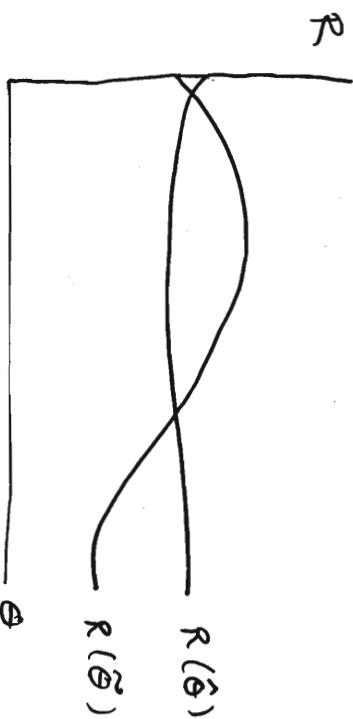
Many common estimators are inadmissible!

e.g. usual OLS estimators of β & σ^2

In linear regression, under quadratic loss, if $K \geq 3$! (Stein, 1956)

(vi) Mini-Max -

Typically, risk funcs. of estimators of interest cross -



which estimator is preferred? We might average the risk (over θ), but what weights?
(Return to this as Bayesians!)

One (conservative) option -

$\hat{\theta}$ is Mini-max within family of estimators, for specific loss fctn. if

$$\text{max. } R(\hat{\theta}) = \min \{ \max_{\theta} R(\theta^k) \},$$
 if θ^* in family of interest.

Some Questions:

- * Are mini-max. estimators admissible?
- * Are admissible estimators mini-max.?
- * Are there any systematic, general, principles for constructing estimators that ensure that the estimators will have "good" properties if the types discussed so far?

We'll be considering 3 general principles -

- (a) Maximum Likelihood Estimation
- (b) Method of Moments Estimation
- (c) Bayesian Estimation

(& the associated testing principles).

As a second-best solution, if it is hard to find principles guaranteed to yield "good" estimators in the above senses, maybe we can find principles that do well for large in

15

16