

Consequently
of the
MGF definition

① Σ rv with MGF $\gamma_{\Sigma}(t)$,
(20)

$$\Sigma = a \bar{\Sigma} + b, \quad (a, b \text{ constants})$$

Then at every value of t for which $\gamma_{\Sigma}(at)$ is finite,

$$\gamma_{\Sigma}(t) = e^{bt} \gamma_{\Sigma}(at). \quad \boxed{\text{Example}}$$

Σ - Binomial (n, p), $\bar{\Sigma} = \sum_{i=1}^n S_i$,

$S_i \stackrel{\text{IID}}{\sim} \text{Bernoulli}(p)$ γ_{S_i} MGF of S_i
($i=1, \dots, n$)

is easy: $\gamma_{S_i}(t) = E(e^{tS_i})$
 $= e^{t \cdot 1} p(S_i=1)$

This way the
law of the
unconscious
statistician

$$+ e^{t \cdot 0} p(S_i=0)$$

$$= [p e^t + (1-p)]$$

② $\underline{X}_1, \dots, \underline{X}_n$ independent w., MGF

of \underline{X}_i is $\mathcal{M}_{\underline{X}_i}(t)$, $\underline{\Sigma} = \sum_{i=1}^n \underline{X}_i$,

MGF of $\underline{\Sigma}$ is $\mathcal{M}_{\underline{\Sigma}}(t) +$ for every
t such that $\mathcal{M}_{\underline{X}_i}(t)$ is finite for all
 $i = 1, \dots, n$.

$$i=1, \dots, n, \quad \mathcal{M}_{\underline{\Sigma}}(t) = \prod_{i=1}^n \mathcal{M}_{\underline{X}_i}(t).$$

~~MGF of~~
Binomial,
continued

Since the S_i are IID,

$$\mathcal{M}_{\underline{\Sigma}}(t) \stackrel{\text{IID}}{=} \prod_{i=1}^n \mathcal{M}_{S_i}(t)$$

Now, as
before, we
just crank out
the derivatives:

$$\stackrel{\text{IID}}{=} \prod_{i=1}^n [pe^t + (1-p)]$$

$$\stackrel{\text{IID}}{=} [pe^t + (1-p)]^n$$

$$E(\bar{X}) = \left(\frac{d}{dt} \mathbb{E}_X(t) \right) \Big|_{t=0} = \frac{d}{dt} \left[p e^t + (1-p) \right] \Big|_{t=0}$$

203

$$E(\bar{X}^2) = \frac{d^2}{dt^2} \left[p e^t + (1-p) \right] \Big|_{t=0} = np [1 + (n-1)p]$$

$$\therefore V(\bar{X}) = E(\bar{X}^2) - [E(\bar{X})]^2$$

$$= np + n(n-1)p^2 - n^2 p^2$$

~~$$= np + \cancel{n^2 p^2} - np^2 - \cancel{n^2 p^2}$$~~

$$= n(p - p^2) = np(1-p)$$

$$E(\bar{X}^3) = \left(\frac{d^3}{dt^3} \left[p e^t + (1-p) \right] \right) \Big|_{t=0} =$$

PF

T-S

(uglier
uglier)

~~$$= np [1 + (n-2)(n-1)p^2 + 3p(n-1)]$$~~

③ \mathbb{X} has MGF $\gamma_{\mathbb{X}}(t)$, finite in an open interval around $t=0$.

\mathbb{Y} has MGF $\gamma_{\mathbb{Y}}(t)$

iff \mathbb{X}, \mathbb{Y} have identical probability distributions

so the MGF (if it exists) uniquely characterizes a random variable.

Mean
variance
median

we've already made some contrasts between the mean and the median of a distribution;

here are 2 more things worth saying.

(CDF $F_{\mathbb{X}}$)

① \mathbb{X} rv with values in an interval I ;
 $h(x)$ $1:1$ function on I , $\Rightarrow I = h(\mathbb{X})$;

if $m_{\bar{X}}$ is ④ median of \bar{X} (ie,

(205)

if $m_{\bar{X}} = F_{\bar{X}}^{-1}(\frac{1}{2})$, then $h(m_{\bar{X}})$ is
a median of $I = h(\bar{X})$. This is

not in general true of the mean,
as we have already seen:

$$E[h(\bar{X})] \neq h[E(\bar{X})]$$

unless $h(x) = ax + b$

Prediction
variable
 \bar{X} rv with
mean $\mu_{\bar{X}}$, SD $\sigma_{\bar{X}}$

Before \bar{X} is observed, suppose your job
is to predict what its value will be;
what should you do? How can you tell
if a prediction is good?

let's say you pick the number \hat{x} (fixed known constant) before \bar{X} is observed.

Then, after \bar{X} owing, your prediction error would be $(\hat{x} - \bar{X})$, which might be either positive or negative. [one]

possible criterion for goodness would be to find \hat{x} such that $E(\hat{x} - \bar{X}) = 0$.

Def] The bias of \hat{x} as a prediction for \bar{X} is $\text{bias}(\hat{x}) \triangleq E(\hat{x} - \bar{X})$.

Def] Your prediction \hat{x} is unbiased

if $\text{bias}(\hat{x}) = 0$. [Clearly, to achieve this just choose $\hat{x} = E(\bar{X})$.]

Another possible criterion for goodness (207)
would be to find \hat{x} such that $E(\hat{x} - \bar{X})^2$

is small. [Def.] $E[(\hat{x} - \bar{X})^2]$ is called the
mean squared error (MSE) of \hat{x} as
a prediction for \bar{X} . [Small MSE means:]

The \hat{x} that minimizes MSE is $\hat{x} = E(\bar{X})$.

Small proof

$$\begin{aligned} E[(\hat{x} - \bar{X})^2] &= E(\hat{x}^2 - 2\hat{x}\bar{X} + \bar{X}^2) \\ &= \hat{x}^2 - 2\hat{x}E(\bar{X}) + E(\bar{X}^2) \end{aligned}$$

This is a quadratic function of \hat{x} :

$$\frac{\partial}{\partial \hat{x}} E[(\hat{x} - \bar{X})^2] = 2\hat{x} - 2E(\bar{X}) = 0$$

iff

$$\hat{x} = E(\bar{X})$$

$$\frac{\partial^2}{\partial \hat{x}^2} = 2 > 0$$

so $E(\bar{X})$ is a minimum

Also easy
to show

$$\text{MSE}(\hat{x}) = E(\hat{x} - \bar{x})^2 \quad (208)$$

$$= V(\bar{x}) + (\text{bias}(\hat{x}))^2$$

So the choice $\hat{x} = E(\bar{x})$ ^{both} minimizes

$\text{MSE}(\hat{x})$ and achieves 0 bias, and

with this choice $\text{MSE}(\hat{x}) = V(\bar{x}) = \sigma_{\bar{x}}^2$

A different criterion for a good prediction \hat{x}
would be to find \hat{x} such
that $E[|\hat{x} - \bar{x}|]$ is small. Definition

$E|\hat{x} - \bar{x}|$ is called the mean absolute error (MAE) of \hat{x} as a prediction for \bar{x}

Another } \bar{X} rv with finite mean $\mu_{\bar{X}}$; 209
 small
 theorem } let $m_{\bar{X}}$ be (a/the) median of \bar{X} .

\rightarrow the \hat{x} that minimizes $MAE(\hat{x})$

is (a/the) median $m_{\bar{X}}$. why
reminder: a/the ?

Careful
 definition
 of median

\bar{X} rv + every number n
 such that

$$P(\bar{X} \leq n) \geq \frac{1}{2} \text{ and } P(\bar{X} \geq n) \geq \frac{1}{2}$$

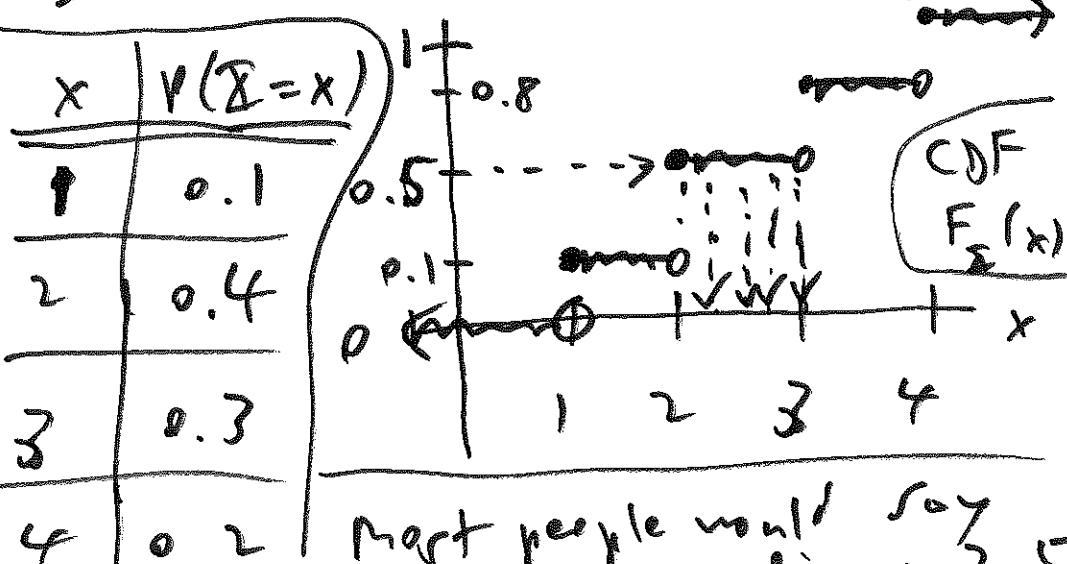
is a median of the dist. of \bar{X}

Example
 of nonunique
 median

All $2 \leq x < 3$

have $F_{\bar{X}}(x) = \frac{1}{2}$

\bar{X} discrete on $\{1, 2, 3, 4\}$



which is better criterion, MSE or MAE?

There is ^{universal} no right answer (210)

to this question: it depends on the real-world consequences of your prediction errors

$(\hat{x} - x)$; quantifying these consequences involves the creation of a utility function,

which we'll examine later.

Covariance & correlation

Independence of 2 or more RVs is a special case of a more general reality, in which (your uncertainty about something) and (your uncertainty about something else) are related. Let's see how to quantify such relationships.

Def. ξ, η rv with finite means μ_ξ and $\mu_\eta = E(\eta)$. The covariance of $E(\xi)$

ξ and η , written $C(\xi, \eta)$, is defined as

If we use
 $\text{Cov}(\xi, \eta) = E[(\xi - \mu_\xi)(\eta - \mu_\eta)]$, as long as this expectation exists

Consequence of this definition

$$\textcircled{1} (\xi - \mu_\xi) \cdot (\eta - \mu_\eta) =$$

$$\xi \cdot \eta - \mu_\xi \cdot \eta - \mu_\eta \cdot \xi + \mu_\xi \mu_\eta$$

$$\therefore C(\xi, \eta) = E(\xi \eta) - \mu_\xi E(\eta) - \mu_\eta E(\xi)$$

$$= E(\xi \eta) - \mu_\xi \mu_\eta - \mu_\xi \mu_\eta + \mu_\xi \mu_\eta$$

$C(\xi, \eta) = E(\xi \eta) - \mu_\xi \mu_\eta$ easier formula
 (expectation of product - product of expectations) to compute with

(2) Sufficient condition for $C(\bar{x}, \bar{y})$ to exist: $\sigma_{\bar{x}}^2 < \infty$ and $\sigma_{\bar{y}}^2 < \infty$. (212)

(3) Covariance

is a good start at measuring strength of relationship, but it has a big flaw: its value depends on the units of measurement of \bar{x} and \bar{y} .

Example: \bar{x} = education level
(years of school ^{fully} completed)

Example:

\bar{y} = yearly income (\$)

\bar{x} = temperature

$C(\bar{x}, \bar{y})$ comes out in

(years) · (\$) ??

in $^{\circ}\text{C}$

\bar{y} = relative humidity (%)

If you change your mind & measure temperature in $^{\circ}\text{F} = \frac{9}{5}^{\circ}\text{C} + 32$,
 $C(\bar{x}, \bar{y}) = C\left(\frac{9}{5}\bar{x} + 32, \bar{y}\right) \neq C(\bar{x}, \bar{y})$

Easy to show that if a, b are fixed constants (213)
 then $C(a\bar{x} + b, \bar{y}) = aC(\bar{x}, \bar{y})$ so
 $C(\bar{x}', \bar{y}) = 1.8 \cdot C(\bar{x}, \bar{y})$, i.e. you can
 make the association
 between temperature & relative
 humidity seem larger just by switching
 from $^{\circ}\text{C}$ to $^{\circ}\text{F}$ (??)

Easy fix:

Def The process of converting a w \bar{x}
 to standard units (say) is achieved with

the linear transformation $\bar{x}' = \frac{\bar{x} - E(\bar{x})}{SD(\bar{x})}$

(or by $\sigma_{\bar{x}} < \infty$, this
 is a meaningful definition)

$$= \frac{\bar{x} - \bar{x}}{\sigma_{\bar{x}}}$$

$$E(\bar{x}') = 0, V(\bar{x}') = 1 = SD(\bar{x}')$$

Def. \bar{X}, \bar{I} rv with finite variance (214)
 $\sigma_{\bar{X}}^2$ and $\sigma_{\bar{I}}^2$ (and therefore finite means
 $\mu_{\bar{X}}$ and $\mu_{\bar{I}}$) \rightarrow the correlation of \bar{X}

and \bar{I} is $\rho(\bar{X}, \bar{I}) = E\left[\left(\frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}}\right) \cdot \left(\frac{\bar{I} - \mu_{\bar{I}}}{\sigma_{\bar{I}}}\right)\right]$

With this definition,
the correlation is

$$= \frac{C(\bar{X}, \bar{I})}{\sigma_{\bar{X}} \cdot \sigma_{\bar{I}}}$$

invariant to linear

transformation of either variable (both):

for any constants $a, c \neq 0$ and b, d ,

$$\rho(a\bar{X} + b, c\bar{I} + d) = \rho(\bar{X}, \bar{I}).$$

(If $a < 0$, $\rho(a\bar{X} + b, \bar{I}) = -\rho(\bar{X}, \bar{I}).$)

Consequences
of the
correlation
definition

① Cauchy-Schwarz inequality:
For all $\mathbf{r} \in \mathbb{R}, \mathbb{I}$ for which
 $E(\mathbf{r}\mathbf{I})$ exists, $(E(\mathbf{r}\mathbf{I}))^2 \leq [E(\mathbf{r})]^2 \cdot [E(\mathbf{I})]^2$.

from which

$$[(C(\mathbf{r}, \mathbf{I}))^2 \leq \sigma_{\mathbf{r}}^2 \cdot \sigma_{\mathbf{I}}^2]$$

Karl Schwarz
(1843-1921)
German mathematician
(associated)

$$\rightarrow -1 \leq \rho(\mathbf{r}, \mathbf{I}) \leq +1.$$

Def $\rho(\mathbf{r}, \mathbf{I}) > 0 \leftrightarrow \mathbf{r}, \mathbf{I}$ positively correlated

$\rho(\mathbf{r}, \mathbf{I}) < 0 \leftrightarrow \mathbf{r}, \mathbf{I}$ negatively correlated

$\rho(\mathbf{r}, \mathbf{I}) = 0 \leftrightarrow \mathbf{r}, \mathbf{I}$ uncorrelated

② \mathbf{r}, \mathbf{I} independent with

$$\begin{aligned} 0 < \sigma_{\mathbf{r}}^2 < \infty \\ 0 < \sigma_{\mathbf{I}}^2 < \infty \end{aligned}$$

$$\rightarrow C(\mathbf{r}, \mathbf{I}) = \rho(\mathbf{r}, \mathbf{I}) = 0$$

so independence implies correlation, ②/6
 but (interestingly) not the converse:

Example: $\underline{X} \sim \text{Uniform}\{-1, 0, +1\}$, $\underline{Y} = \underline{X}^2$
 $E(\underline{X}) = 0$

$\rightarrow \underline{X}, \underline{Y}$ clearly dependent since \underline{X} completely determines \underline{Y} , but $E(\underline{X}\underline{Y}) = E(\underline{X}^3)$

(since \underline{X} and \underline{X}^3 are identically distributed) and thus

$$= E(\underline{X}) = 0$$

$$C(\underline{X}, \underline{Y}) = \underbrace{E(\underline{X}\underline{Y})}_{0} - \underbrace{E(\underline{X}) \cdot E(\underline{Y})}_{0} = 0$$

so $\rho(\underline{X}, \underline{Y}) = \frac{C(\underline{X}, \underline{Y})}{\sigma_{\underline{X}} \sigma_{\underline{Y}}} = 0$ and $\underline{X}, \underline{Y}$ are uncorrelated.

③ \underline{X} rv with $0 < \sigma_{\underline{X}}^2 < \infty$, $\underline{Y} = a\underline{X} + b$
 for $\{a \neq 0\}$ constants $\rightarrow (a > 0) \rho(\underline{X}, \underline{Y}) = +1$

$$(a < 0) \rho(X, Y) = -1 \quad \text{so } \rho(Y, X) \quad (21)$$

measures the strength of linear association between X and Y .

Important:

$$X, Y \text{ rv}, \sigma_X^2 < \infty, \sigma_Y^2 < \infty \rightarrow$$

$$V(X+Y) = V(X) + V(Y) + 2C(X, Y)$$

$$\textcircled{5} \left(\begin{array}{c} a, b, c \\ \text{any constants} \end{array} \right) C(aX + bY) = abC(X, Y)$$

$$\sigma_X^2 < \infty, \sigma_Y^2 < \infty \rightarrow V(aX + bY + c) =$$

Special case: $a^2 V(X) + b^2 V(Y) + 2abC(X, Y)$

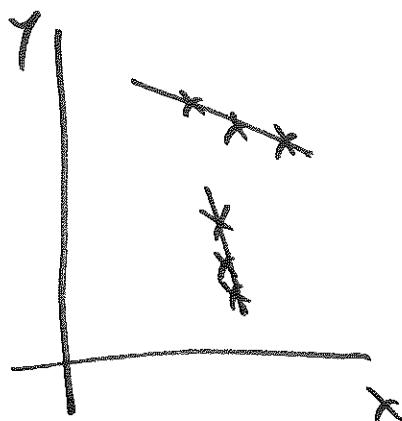
$$V(X-Y) = V(X) + V(Y) - 2C(X, Y).$$

⑥ $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that $(\mathbf{x}_i, \mathbf{x}_j)$ uncorrelated (218)

$$\text{for all } 1 \leq i \neq j \leq n \Rightarrow \sqrt{\left(\sum_{i=1}^n \mathbf{x}_{ii}\right)} = \sum_{i=1}^n \sqrt{\mathbf{x}_{ii}}$$

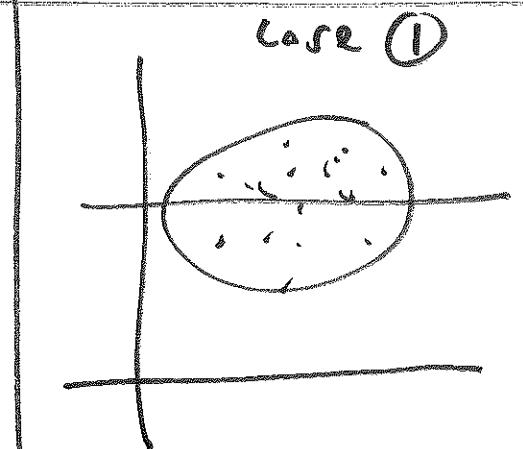
⑦

$$P(\mathbf{x}, \mathbf{x}) = -1$$

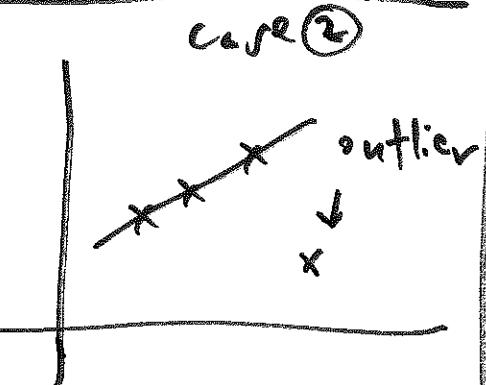


points in
scatterplot
sample from
 $f_{\mathbf{x}, \mathbf{x}}(\mathbf{x}, \mathbf{y})$
all fall on line
with negative
slope (not
necessarily
 -1)

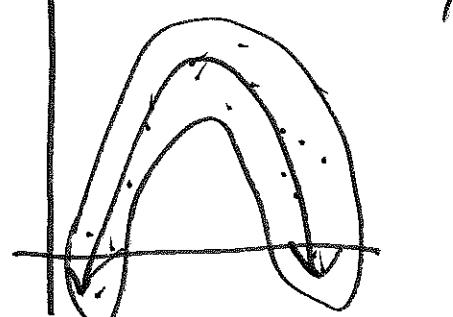
$$P(\mathbf{x}, \mathbf{x}) = 0$$



case ①

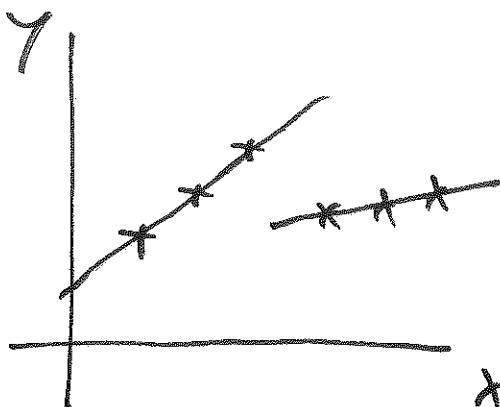


case ②



non-linearity

$$P(\mathbf{x}, \mathbf{x}) = +1$$



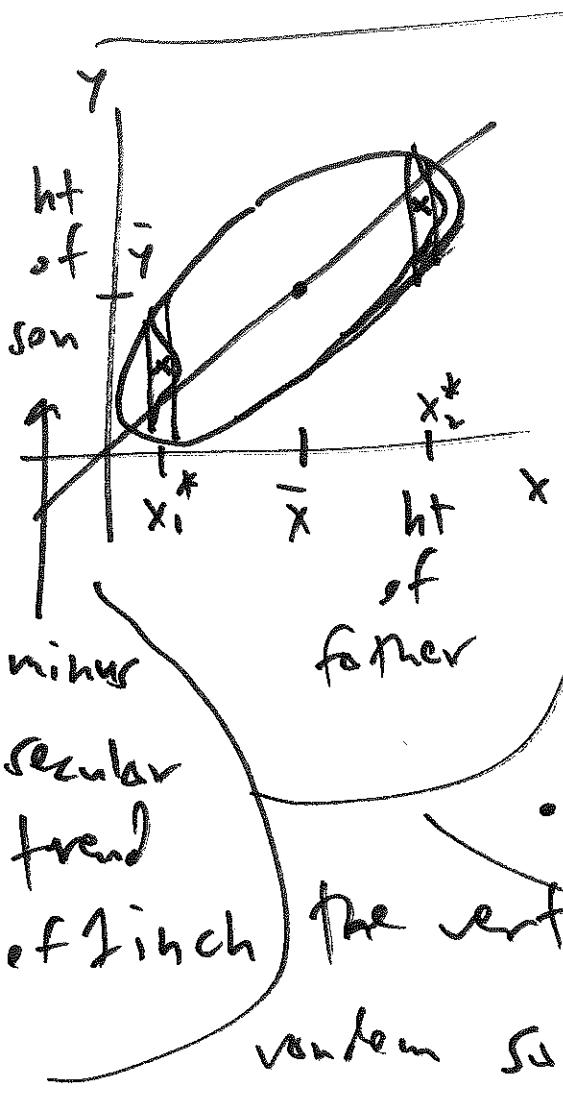
points in
scatterplot
sample from
 $f_{\mathbf{x}, \mathbf{x}}(\mathbf{x}, \mathbf{y})$

all fall on line
with positive
slope (not
necessarily
 $+1$)

Conditional Expectation

X, Y related rvs (not independent). Then there is information in X for predicting Y ; i.e., we should be able to find some function $\lambda: \mathbb{R} \rightarrow \mathbb{R}$ such that $\lambda(X)$ is "close" in some sense to Y — what is the optimal λ ?

(219)



Galton example again:

Galton divided the elliptical scatterplot up into a bunch of vertical strips, e.g., the one over x_1^* or the other one over x_2^* . ~~Observe~~ The points in the vertical strip over x_2^* are a random sample from the conditions!

220

distribution of Σ given $\bar{X} = \bar{x}_2^*$, $f_{\Sigma|\bar{X}}(\gamma | \bar{x} = \bar{x}^*)$

Galton knew about the small theorem

(text on p. 207): the number \hat{w} that minimizes

(MSE) the mean squares error $E[(\hat{w} - \bar{w})^2]$ of \hat{w}

" w " prediction for \bar{w} is $\hat{w} = E(\bar{w})$.

So he adopted MSE as his measure of "close"

and concluded that the $\hat{\gamma}$ that minimizes

the MSE $E[(\hat{\gamma} - \Sigma)^2]$ in the vertical strip

defined by $x = \bar{x}_2^*$ must be the conditional

mean, or conditional expectation, of the

$\sim (\Sigma | \bar{X} = \bar{x}_2^*)$ Def. $\bar{x}, \Sigma \sim n, \Sigma$ finite mean +

$\left\{ \begin{array}{l} \text{conditional expectation} \\ (\text{mean}) \text{ of } \Sigma \text{ given } \bar{X} = \bar{x} \end{array} \right\} = E(\Sigma | \bar{x})$ is just

the expectation of the conditional distribution,

$f_{\Sigma|\Xi}(y|x)$ of Σ given $\Xi = x$,

namely $E(\Sigma|x) = \int_{\mathbb{R}} y f_{\Sigma|\Xi}(y|x) dy$

for continuous $(\Sigma|\Xi=x)$

and $E(\Sigma|x) = \sum_{\text{all } y} y f_{\Sigma|\Xi}(y|x)$

for discrete $(\Sigma|\Xi=x)$

so far, $E(\Sigma|x)$ is just a constant,
equal to the conditional mean of Σ

when Ξ is x Def. $h(x) \stackrel{\Delta}{=} E(\Sigma|\Xi=x)$

then the w $E(\Sigma|\Xi) \stackrel{\Delta}{=} h(\Xi)$ is the
conditional expectation of Σ given Ξ

Clinical trial

example,

continued

$(n_c + n_T)$ people^(a) who are similar in all relevant ways to population P = {all adult patients with disease A}

and (b) who consent to participate in your clinical trial are randomized, n_c to ^{the} control group and n_T to ^{the} treatment group.

outcome of interest is dichotomous:

let θ be the proportion of successes you would have seen if you could have put (everybody in P) into your treatment group; θ is unknown.

(success)	$I =$ disease went into remission
(failure)	$\bar{I} =$ didn't

let $S_i = \begin{cases} 1 & \text{if patient } i \text{ is in the actual } T \text{ group} \\ & \text{had a success} \\ 0 & \text{otherwise} \end{cases}$

then the rvs $(S_i | \theta)$ are IID Bernoulli(θ)⁽²³⁾

and the rv $S = \sum_{i=1}^{n_T} S_i$ has a conditional

binomial dist: $(S | \theta) \sim \text{Binomial}(n_T, \theta)$

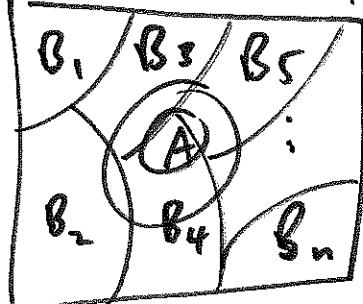
It's meaningful to talk about the conditional expectation rv. $E(S | \theta) = n_T \theta$ (a linear function of θ),

and - via Bayes' Theorem - it's even more meaningful to talk about the conditional expectation rv. $E(\theta | S)$ (more about this later)

and the constant $E(\theta | S = s)$. / Important

(5 Aug)

Remember the law of total prob.!



$$P(A) = \sum_{i=1}^n P(B_i) P(A|B_i)$$

(LTP)

Consequence
of the
def. of
conditional
expectation