

③ X has MGF $\psi_X(t)$, finite in an open interval around $t=0$.
 Y has MGF $\psi_Y(t)$.

then $\psi_X(t) = \psi_Y(t) \iff X, Y$ have identical probability distributions

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

Mean
 versus
 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_X)
 ① X rv with values in an interval I ;
 $h(x)$ 1-1 function on I , $h(I) = h(X)$;

if m_X is a median of X (ie, (205)

if $m_X = F_X^{-1}(\frac{1}{2})$, then $h(m_X)$ is
a median of $Y = h(X)$. This is

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

$$E[h(X)] \neq h(E(X))$$

unless $h(x) = ax + b$

Prediction
~~Definition~~
 X rv with
mean μ_X , SD σ_X

Before 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} ²⁰⁶ \leftarrow x -hat
(a fixed known constant) before X is observed.

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

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

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

Def) Your prediction \hat{x} is unbiased

if $\text{bias}(\hat{x}) = 0$.

Clearly, to achieve this just choose $\hat{x} = E(X)$.

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

is small. (Gauss) Def. $E[(\hat{x} - X)^2]$ is called the

mean/squared error (MSE) of \hat{x} as

a prediction for X . Small ~~the~~ theorem:

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

Small proof

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

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

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

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

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

so $E(X)$ is a minimum

Also easy to show

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

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

So the choice $\hat{x} = E(X)$ ^{both} minimizes $MSE(\hat{x})$ and achieves 0 bias, and with this choice $MSE(\hat{x}) = V(X) = \sigma_X^2$

A different criterion

Yet another possible criterion for a good prediction \hat{x} would be to find \hat{x} such

that $E[|\hat{x} - X|]$ is small.

(Laplace)

Definition

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

Another small theorem

X rv with finite mean μ_X ; (209)
 let m_X be (a/the) median of X ;

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

is (a/the) median m_X . Reminder: why a/the?

Careful definition of median

X rv \rightarrow every number m such that

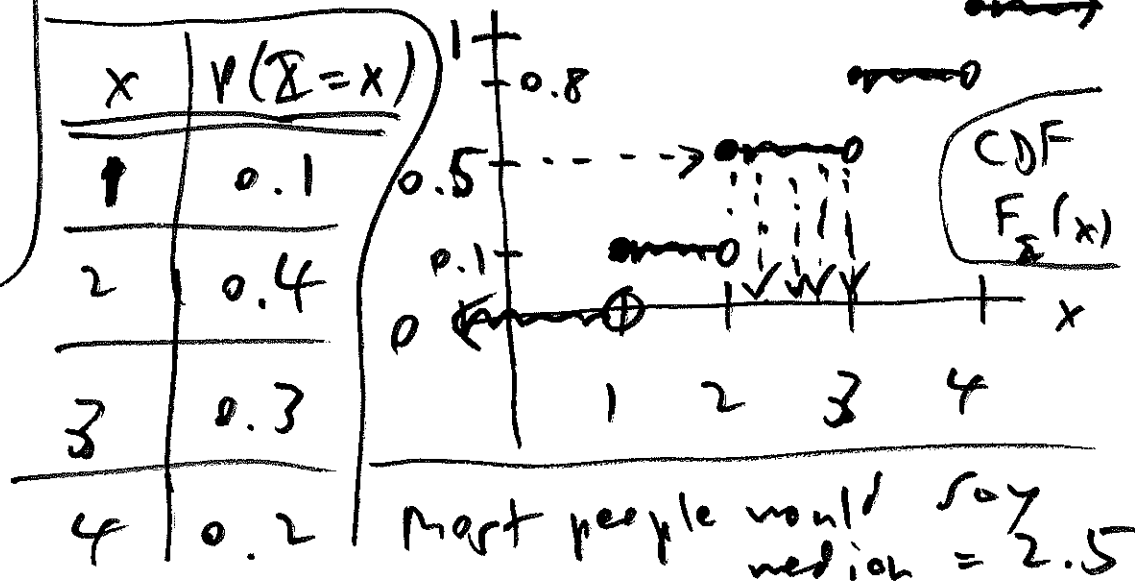
$$P(X \leq m) \geq \frac{1}{2} \text{ and } P(X \geq m) \geq \frac{1}{2}$$

is a median of the dist. of X

Example of nonunique median

All $2 \leq x < 3$ have $F_X(x) = \frac{1}{2}$

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



Most people would say median = 2.5

which is
a 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 ^{briefly} 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. X, Y rv with finite means μ_X and $\mu_Y = E(Y)$. The covariance of X and Y , written $C(X, Y)$, is defined as

If use $Cov(X, Y)$ $C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$, as long as this expectation exists

Consequences of this definition

① $(X - \mu_X) \cdot (Y - \mu_Y) = X \cdot Y - \mu_X \cdot Y - \mu_Y \cdot X + \mu_X \mu_Y$

so $C(X, Y) = E(XY) - \mu_X E(Y) - \mu_Y E(X) + \mu_X \mu_Y$

$C(X, Y) = E(XY) - \mu_X \mu_Y$ (expectation of product - product of expectations) easier formula to compute with

② Sufficient condition for $C(X, Y)$ to (212)

exist: $\sigma_X^2 < \infty$ and $\sigma_Y^2 < \infty$.

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 X and Y

Example:

$X = \overset{\text{max}}{\text{temperature}}$

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

$Y = \overset{\text{max}}{\text{humidity}} (\%)$

Example: $X = \text{education level}$
(years of schooling completed)

$Y = \text{yearly income } (\$)$

$C(X, Y)$ comes out in

(years) \cdot (\$) (??)

If you change your mind & measure temperature X' in $^{\circ}\text{F} = \frac{9}{5}C + 32$,

$$C(X', Y) = C\left(\frac{9}{5}X + 32, Y\right) \neq C(X, Y)$$

Easy to show that if a, d are ^{fixed} constants (23)

then $C(aX + b, Y) = a C(X, Y)$ so

$$C(X', Y) = 1.8 \cdot C(X, Y), \text{ i.e. you can}$$

of $^{\circ}F$ \nearrow $^{\circ}C$ make the association between temperature & relative humidity seem larger just by switching from $^{\circ}C$ to $^{\circ}F$ (???)

Easy fix:

Def The process of converting a rv X to standard units (SU) is achieved with

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

$SD(X)$

(as long as $\sigma_X < \infty$, this is a meaningful definition)

$$= \frac{X - \mu_X}{\sigma_X}$$

$$E(X') = 0, \quad V(X') = 1 = SD(X')$$

Def. / X, Y rv with finite variances (214)
 σ_X^2 and σ_Y^2 (and therefore finite means
 μ_X and μ_Y) \rightarrow the correlation of X

and Y is $\rho(X, Y) = E \left[\left(\frac{X - \mu_X}{\sigma_X} \right) \cdot \left(\frac{Y - \mu_Y}{\sigma_Y} \right) \right]$

With this definition,
the correlation is
invariant to linear

$$= \frac{C(X, Y)}{\sigma_X \cdot \sigma_Y}$$

transformation of either variable (both):

for any constants $a, c > 0$ and $b, d,$

$$\rho(aX + b, cY + d) = \rho(X, Y).$$

(If $a < 0$, $\rho(aX + b, Y) = -\rho(X, Y)$.)

Consequences
of the
correlation
definition

① Cauchy-Schwarz inequality (215)

For all $r.v. X, Y$ for which

$$E(XY) \text{ exists, } (E(XY))^2 \leq [E(X)]^2 \cdot [E(Y)]^2$$

from which $[C(X, Y)]^2 \leq \sigma_X^2 \cdot \sigma_Y^2$

and $-1 \leq \rho(X, Y) \leq +1$

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

Def $\rho(X, Y) > 0 \leftrightarrow X, Y$ positively
correlated

$\rho(X, Y) < 0 \leftrightarrow X, Y$ negatively
correlated

$\rho(X, Y) = 0 \leftrightarrow X, Y$ uncorrelated

② X, Y independent $r.v.$ with $\left\{ \begin{array}{l} 0 < \sigma_X^2 < \infty \\ 0 < \sigma_Y^2 < \infty \end{array} \right\}$

$\rightarrow C(X, Y) = \rho(X, Y) = 0$

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

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

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

(since X and X^3 are
identically distributed) $= E(X) = 0$
and thus

$$C(X, Y) = \underbrace{E(XY)}_0 - \underbrace{E(X)}_0 \cdot E(Y) = 0$$

$$\therefore \rho(X, Y) = \frac{C(X, Y)}{\sigma_X \sigma_Y} = 0 \quad \text{and } X, Y \text{ are uncorrelated!}$$

③ $X \sim N$ with $0 < \sigma_X^2 < \infty$, $Y = aX + b$
for $\begin{cases} a \neq 0 \\ b \end{cases}$ constants $\rightarrow (a > 0) \rho(X, Y) = +1$

$$(a < \infty) \rho(X, Y) = -1 \quad \text{so} \quad \rho(X, Y) \quad (217)$$

measures the strength of linear association between X and Y .

(4) Important:

(if)

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

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

(5) $\begin{matrix} a, b, c \\ \text{any} \\ \text{constants} \end{matrix}$ $C(aX, bY) = ab C(X, Y)$

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

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

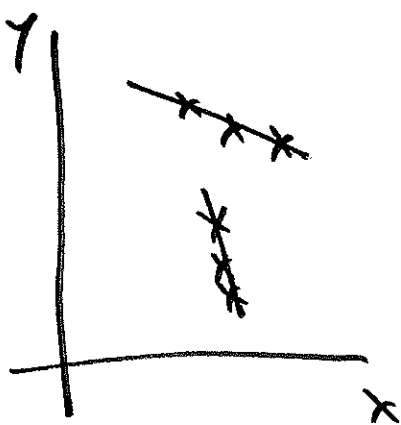
$$V(X-Y) = V(X) + V(Y) - 2C(X, Y)$$

⑥ (Def) X_1, \dots, X_n such that (X_i, X_j) uncorrelated (218)

for all $1 \leq i \neq j \leq n \rightarrow$ (then) $V(\sum_{i=1}^n X_i) = \sum_{i=1}^n V(X_i)$

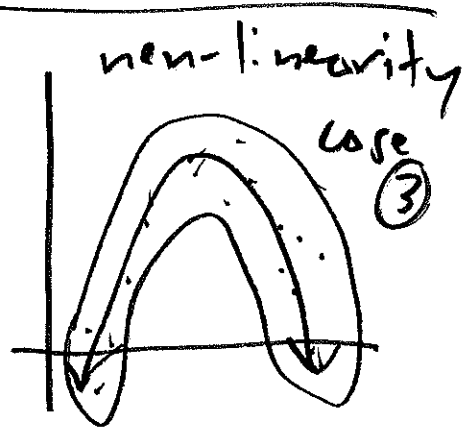
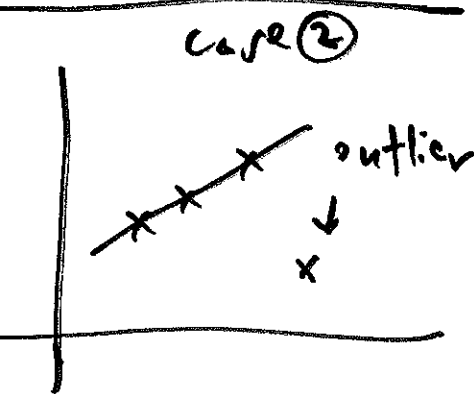
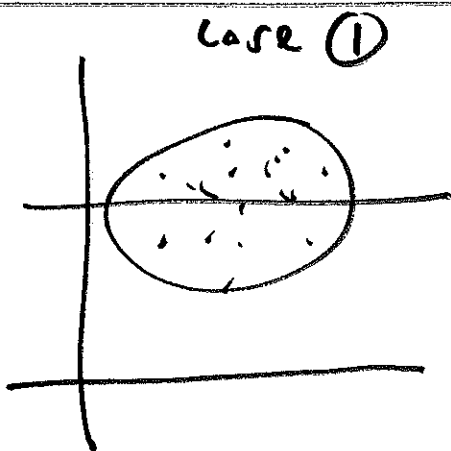
⑦

$\rho(X, Y) = -1$

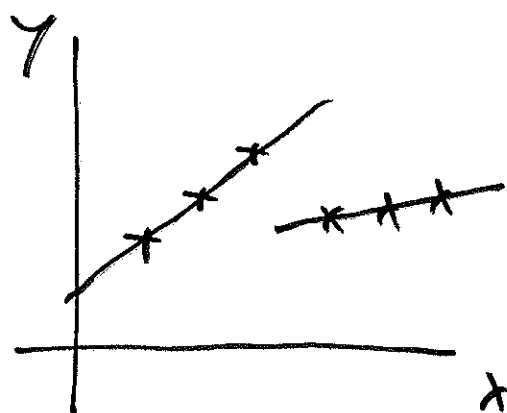


points in scatter plot sample from $f_{X,Y}(x,y)$ all fall on line with negative slope (not necessarily -1)

$\rho(X, Y) = 0$



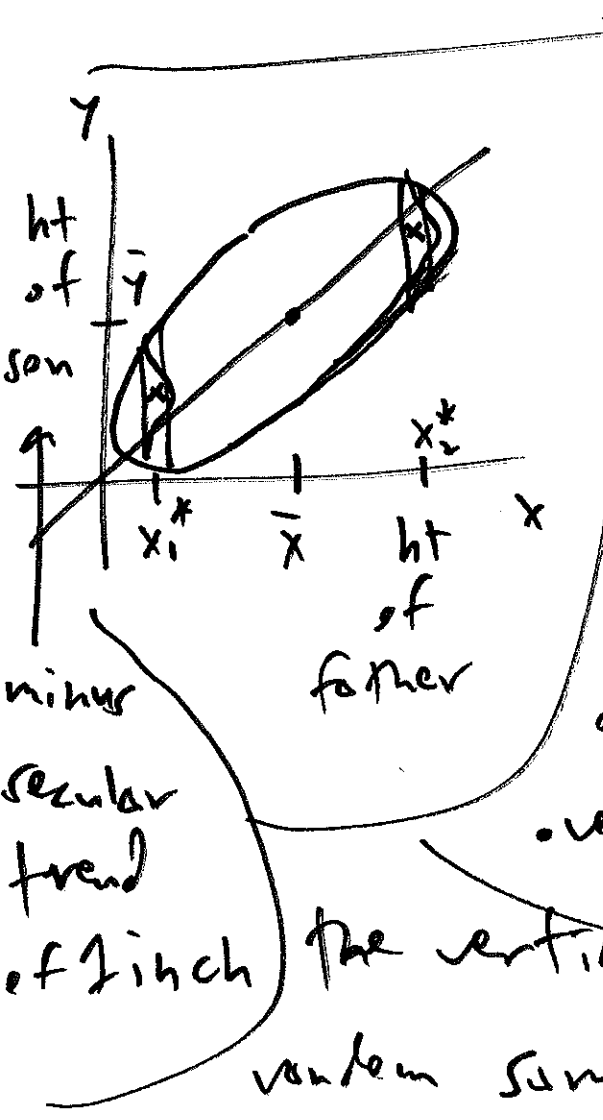
$\rho(X, Y) = +1$



points in scatter plot sample from $f_{X,Y}(x,y)$ all fall on line with positive slope (not necessarily +1)

(21 Aug 17)
 Conditional
 Expectation

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



Galton example ~~graph~~:

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^* .

The points in the vertical strip over x_2^* are a random sample from the conditional

distribution of Y given $X = x_2^*$, $f_{Y|X}(y|x=x_2^*)$ (220)

Galton knew about the small theorem

but on p. (207): the number \hat{w} that minimizes the mean squared error (MSE) $E[(\hat{w} - W)^2]$ of \hat{w} as a prediction for W is $\hat{w} = E(W)$.

So he adopted MSE as his measure of "closeness" and concluded that the \hat{y} that minimizes the MSE $E[(\hat{y} - Y)^2]$ in the vertical strip defined by $x = x_2^*$ must be the conditional mean, or conditional expectation, of the

$v(Y|X = x_2^*)$ Def. $E, \mathbb{R}^n, \mathbb{R}$ finite mean \rightarrow

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

the expectation of the conditional distribution (221)

$f_{Y|X}(y|x)$ of Y given $X=x$,

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

for continuous $(Y|X=x)$

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

for discrete $(Y|X=x)$

So far, $E(Y|x)$ is just a constant,
equal to the conditional mean of Y

the constant

when X is x . Def. $h(x) \triangleq E(Y|X=x)$

then the rv $E(Y|X) \triangleq h(X)$ is the
conditional expectation of Y given X .

Clinical trial example, continued

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

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. (c)

outcome of interest is dichotomous:

(success)	$1 =$ disease went into remission
(failure)	$0 =$ did not

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.

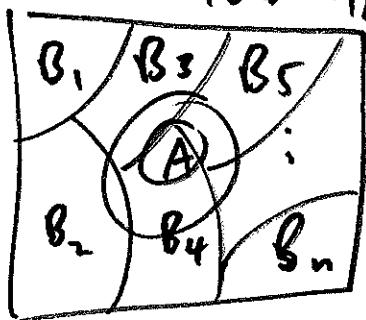
let $S_i = \begin{cases} 1 & \text{if patient } i \text{ is in the actual } \textcircled{T} \text{ group had a success} \\ 0 & \text{otherwise} \end{cases}$

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 $E(\theta | S = s)$.

Remember the Law of Total Prob.!



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

(LTP)

Important
 consequence
 of the
 def. of
 conditional
 expectation

Continuous version of LTP

X, Y continuous r.v. (224)

for which all named densities exist \rightarrow

$$f_Y(y) = \int_{-\infty}^{\infty} f_X(x) \cdot f_{Y|X}(y|x) dx$$

Earlier we agreed that, by definition,

$$E(Y|x) = \int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy$$

So watch the following slightly magical calculation:

$$\begin{aligned} E(Y) &= \int_{-\infty}^{\infty} y f_Y(y) dy \\ &= \int_{-\infty}^{\infty} y \left[\int_{-\infty}^{\infty} f_X(x) f_{Y|X}(y|x) dx \right] dy \end{aligned}$$

ifok to interchange order of integration

$$= \int_{-\infty}^{\infty} f_X(x) \left[\int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy \right] dx$$

$$= \int_{-\infty}^{\infty} f_X(x) \cdot E(Z|x) dx, \text{ and this} \quad (*)$$

is of the form { weighted average of $E(Z|x)$,
with $f_X(x)$ as the weights }

Recall that
continuous
for any r.v. W ,

$$E(W) = \int_{-\infty}^{\infty} w f_W(w) dw$$

and

$$E(h(W)) = \int_{-\infty}^{\infty} h(w) f_W(w) dw \quad (\text{LOTUS})$$

so $(*)$ is just

$$E_X [E(Z|X)]$$

and we have shown that (Adam)

$$E(Z) = E_X [E(Z|X)]$$

This is referred to as part **(1)** of the
double expectation theorem; strangely, I
don't even mention that name, calling it instead
the LTP for expectations.

I need to postpone examples of these conditional expectation calculations until we've covered more standard distributions.

~~Def~~ X, Y r.v. such that $f_{Y|X}(y|x)$ exists \rightarrow it makes sense to speak not only of $E(Y|x)$, the mean of $f_{Y|X}(y|x)$, but also of the variance of that dist.

Def $V(\overbrace{Y|X}^{\text{the number}}) \triangleq E_X \left\{ [Y - E(Y|x)]^2 | x \right\}$
is called the conditional variance of Y given $X = x$, and the r.v. $V(Y|X)$ is just $g(X)$, the conditional variance of Y given X .

The payoff (formalizing Galton's intuition) (227)

from all of this

Theorem X, Y related r.v.;
want to use some function

$\hat{Y} = d(X)$ to predict Y from X \rightarrow

the prediction $\hat{Y} = d(X)$ that minimizes

the MSE $E(Y - \hat{Y})^2 = E\left\{\left[Y - d(X)\right]^2\right\}$

is $\hat{Y} = d(X) = E(Y|X)$, the conditional expectation of Y given X .

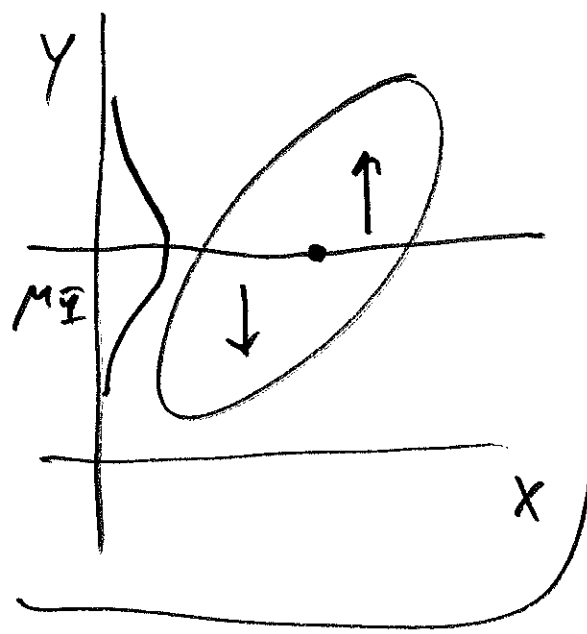
X, Y r.v. such that all of the following expressions exist, \rightarrow

$$V(Y) = E_X[V(Y|X)] + V_X[E(Y|X)].$$

Part (2) of the double expectation theorem

(Eve)

Imagine a 2-part game!



Stage 1 Predict \underline{Y} without knowing \underline{X} . Well, if you buty into MSE as your

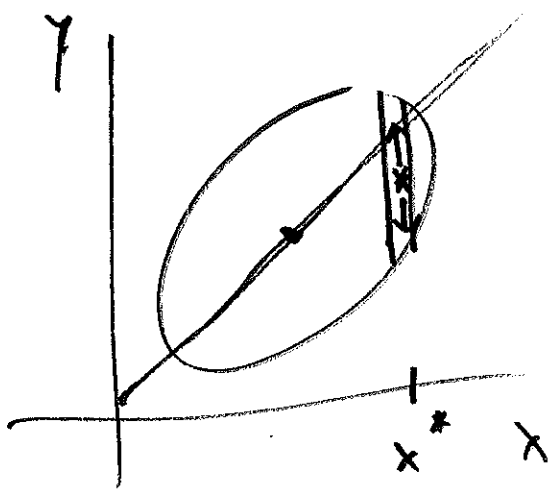
measure of "goodness" of a prediction, we know that you should predict $\hat{\underline{Y}}_{\text{no } \underline{X}} = \mu_{\underline{Y}} = E(\underline{Y})$

and your resulting MSE will be

$$E[(\underline{Y} - \mu_{\underline{Y}})^2] = V(\underline{Y}) = \sigma^2_{\underline{Y}}$$

Stage 2

observe \underline{X} , now predict \underline{Y}



let's say $\underline{X} = x^*$

Then we

know the MSE-optimal

prediction is $\hat{\underline{Y}}_{\underline{X}=x^*} = E(\underline{Y} | \underline{X}=x^*)$

and your resulting MSE will be

$$E \left\{ \left[Y - E(Y|X=x^*) \right]^2 \right\} = \underbrace{V(Y|X=x^*)}_{**}$$

From the vantage point of someone thinking about stage 2 before it happens, X is not yet known, so the expected value of $**$,

namely $E_X [V(Y|X)]$, is the best you can do to guess at how good the stage 2 prediction will be.

The second part of

the double expectation theorem says

$$\underbrace{V(Y)}_{\substack{\uparrow \\ \text{MSE of} \\ \hat{Y}_{no X}}} = \underbrace{E_X [V(Y|X)]}_{\substack{\text{"E(MSE)" of} \\ \hat{Y}_X = E(Y|X)}} + \underbrace{V_X [E(Y|X)]}$$

But since variances are always non-negative,

$$V_X [E(Y|X)] \geq 0, \text{ so}$$

$$E_X [V(Y|X)] + V_X [E(Y|X)] \geq E_X [V(Y|X)]$$

$$V(Y)$$

\geq

"E(MSE)"
of \hat{Y}_X

MSE of $\hat{Y}_{no X}$

Thus you always expect your predictive accuracy to get better (or at least stay the same) when you use $E(Y|X)$ to predict Y .

Another complete switch in subject!

Utility

Q: How to take action sensibly when the consequences are uncertain?