

(2)  $X_1, \dots, X_n$  independent rv, MGF

of  $X_i$  is  $\psi_{X_i}(t)$ ,  $Y = \sum_{i=1}^n X_i$ ,

MGF of  $Y$  is  $\psi_Y(t) \rightarrow$  for every  $t$  such that  $\psi_{X_i}(t)$  is finite for all

$i=1, \dots, n$ ,  $\psi_Y(t) = \prod_{i=1}^n \psi_{X_i}(t)$ .

(18 Aug 17)

MGF of Binomial, continued

Since the  $S_i$  are IID,

$\psi_Y(t) \stackrel{\text{IID}}{=} \prod_{i=1}^n \psi_{S_i}(t)$

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

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

Now, as before, we just crank out the derivative.

$$E(X) = \left( \frac{d}{dt} \psi_X(t) \right) \Big|_{t=0} = \frac{d}{dt} [pe^t + (1-p)]^n \Big|_{t=0} \quad \text{203}$$

$$= np \quad \checkmark$$

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

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

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

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

$$= n(p - p^2) = np(1-p) \quad \checkmark$$

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



(unlike  
& unlike)

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

∴

③  $X$  has MGF  $\psi_X(t)$

finite in an open interval around  $t=0$

$Y$  has MGF  $\psi_Y(t)$

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

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

Mean  
vs  
Median

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

here are 2 more things worth saying.

- (a)  $F_X$
- ①  $X$  rv with values in an interval  $I$ ;  
 $h(x)$  1-1 function on  $I$ ,  $Y = h(X)$ ;

if  $m_{\mathcal{X}}$  is a median of  $\mathcal{X}$  (ie, (205)

if  $m_{\mathcal{X}} = F_{\mathcal{X}}^{-1}(\frac{1}{2})$ , then  $h(m_{\mathcal{X}})$  is

a median of  $\mathcal{Y} = h(\mathcal{X})$ . This is

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

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

unless  $h(x) = ax + b$

$\mathcal{X}$  rv with  
mean  $\mu_{\mathcal{X}}$ , SD  $\sigma_{\mathcal{X}}$

Prediction  
~~prediction~~

Before  $\mathcal{X}$  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  $\hat{x}$  <sup>(206)</sup>  <sub>$\leftarrow x\text{-hat}$</sub>  (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.

---

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 theorem:

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

Small proof

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

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

$$\begin{aligned} \text{MSE}(\hat{x}) &= E(\hat{x} - X)^2 \quad (208) \\ &= V(X) + (\text{bias}(\hat{x}))^2 \end{aligned}$$

So the choice  $\hat{x} = E(X)$  <sup>both</sup> minimize,  $\text{MSE}(\hat{x})$  and achieves 0 bias, and

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

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

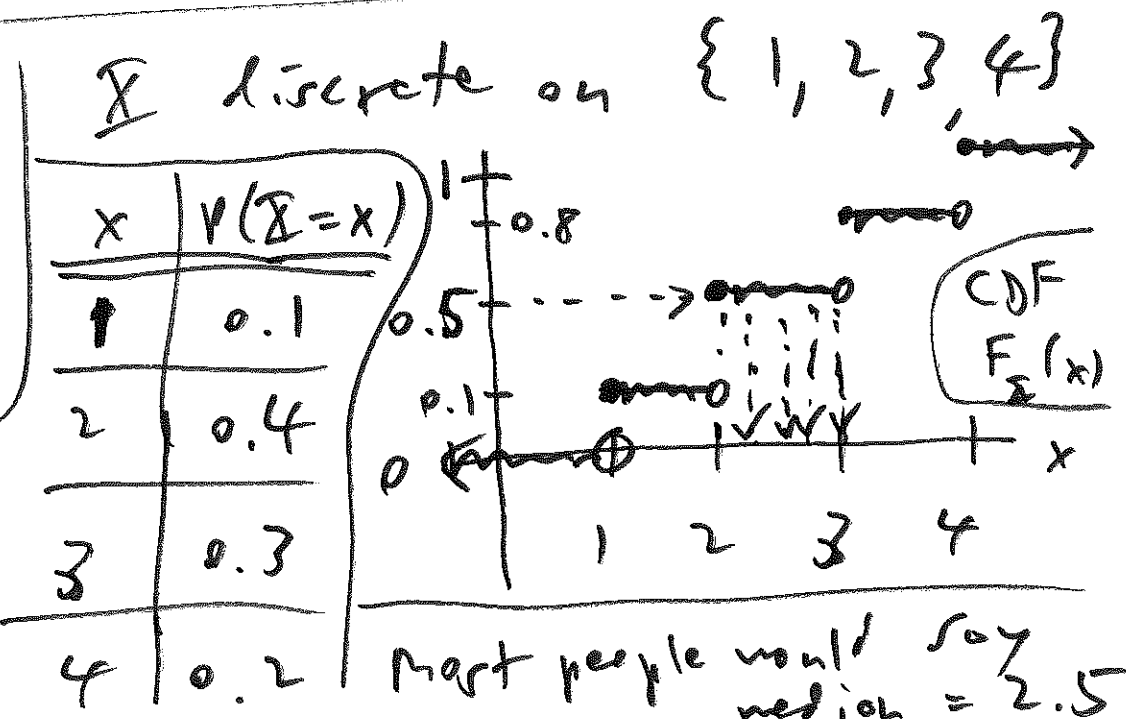
$E|\hat{x} - X|$  is called the mean absolute error (MAE) of  $\hat{x}$  as a prediction for  $X$

Another small theorem )  $X$  rv with finite mean  $\mu_X$ ; (209)  
 let  $m_X$  be (a/the) median of  $X$ ;

$\rightarrow$  the  $\hat{x}$  that minimizes  $M\hat{A}E(\hat{x})$   
 is (a/the) median  $m_X$ . Reminder: why a/the?

Careful definition of median  
 $X$  rv  $\rightarrow$  every number  $m$  such that  
 $P(X \leq m) \geq \frac{1}{2}$  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  $F_X(x) = \frac{1}{2}$





which is  
a better  
criterion,  
MSE or  
MAE?

There is <sup>universal</sup> 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 <sup>briefly</sup> 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  $\mu_X$

and  $\mu_Y = E(Y)$ .

The covariance of  $E(X)$

$X$  and  $Y$ , written  $C(X, Y)$ , is defined as

If we  
cov(X, Y)

$$C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)], \text{ as long as this expectation exists}$$

Consequences  
of this  
definition

$$\textcircled{1} (X - \mu_X) \cdot (Y - \mu_Y) =$$

$$X \cdot Y - \mu_X \cdot Y - \mu_Y \cdot X + \mu_X \mu_Y$$

$$\text{so } C(X, Y) = E(XY) - \mu_X E(Y) - \mu_Y E(X)$$

$$= E(XY) - \mu_X \mu_Y - \mu_Y \mu_X + \mu_X \mu_Y$$

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

much  
easier formula  
to compute  
with

(expectation of product -  
product of expectations)

② Sufficient condition for  $C(X, Y)$  to exist:  $\sigma_X^2 < \infty$  and  $\sigma_Y^2 < \infty$ . (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  $X$  and  $Y$ .

Example:  
 $X = \overset{\text{max}}{\text{temperature}}$   
in  $^{\circ}\text{C}$   
 $Y = \overset{\text{max}}{\text{relative humidity}}$  (%)

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  $^{\circ}\text{F} = \frac{9}{5}C + 32$ ,  
 $C(X', Y) = C(\frac{9}{5}X + 32, Y) \neq C(X, Y)$

Easy to show that if  $a, b$  are <sup>fixed</sup> 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  $^{\circ}F$   $\nearrow$   $\nearrow$   $^{\circ}C$  make the association between temperature & relative humidity seen layer just by switching from  $^{\circ}C$  to  $^{\circ}F$  (???)

Easy fix:

Def The process of converting a rv  $X$  to standard units (54) is achieved with

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

(or by 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  $\sigma_X^2$  and  $\sigma_Y^2$  (and therefore finite means  $\mu_X$  and  $\mu_Y$ )  $\rightarrow$  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 rv  $X, Y$  for which

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

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

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

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

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

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

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

②  $X, Y$  independent rv 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 0 correlation, (2/16)  
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

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

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

---

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

$$(a < 0) \rho(X, Y) = -1 \quad \text{so } \rho(X, Y) \quad (217)$$

measures the strength of linear association between  $X$  and  $Y$ .

④ Important:

⑤ 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  $\begin{matrix} a, b, c \\ \text{any} \\ \text{constants} \end{matrix}$

$$C(aX, bY) = ab C(X, Y)$$

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

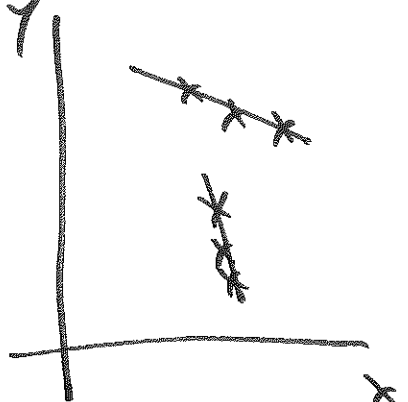
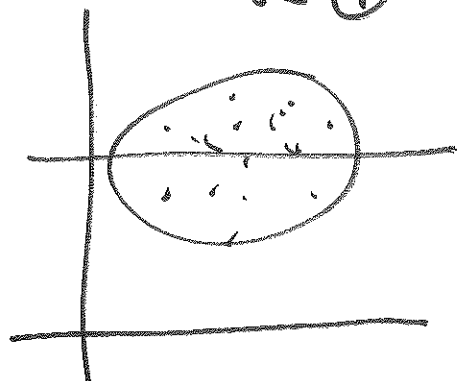
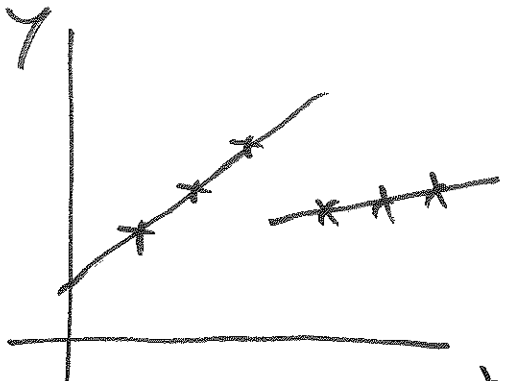
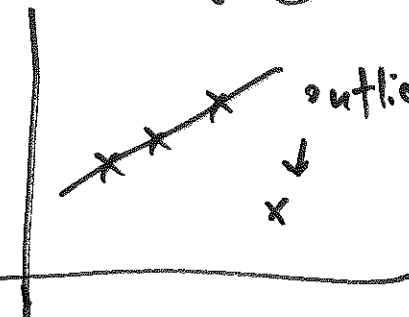
Special case:

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

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



⑥ <sup>(218)</sup>  $X_1, \dots, X_n$  such that  $(X_i, X_j)$  uncorrelated  
 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 $\rho(X, Y) = -1$	$\rho(X, Y) = 0$	$\rho(X, Y) = +1$
	<p>Case ①</p> 	
<p>points in scatterplot sample from <math>f_{X,Y}(x,y)</math> all fall on line with negative slope (not necessarily -1)</p>	<p>Case ②</p> 	<p>points in scatterplot sample from <math>f_{X,Y}(x,y)</math> all fall on line with positive slope (not necessarily +1)</p>
<p>non-linearity</p> 