

Consequences  
of the  
MGF definition

①  $X$  rv with MGF  $\psi_X(t)$ , (20)

$$Y = aX + b, \quad (a, b \text{ constants})$$

Then at every value of  $t$  for which  $\psi_X(at)$  is finite,

$$\psi_Y(t) = e^{bt} \psi_X(at).$$

Example

$X \sim \text{Binomial}(n, p)$ ,  $X = \sum_{i=1}^n S_i$ ,

$S_i \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$   
( $i=1, \dots, n$ )

MGF of  $S_i$

is easy:  $\psi_{S_i}(t) = E(e^{tS_i})$

$$= e^{t \cdot 1} \cdot p(S_i=1)$$

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

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

This was the  
Law of the  
unconscious  
Statistician  
(Lotus)

②  $X_1, \dots, X_n$  independent r.v., 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

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

MGF of Binomial, continued

(18 Aug 17)

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.

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


---


$$= np \checkmark$$

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


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

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

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

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

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



(using  
& y/lie)

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


---

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 (21 May 22)

|                          |  |
|--------------------------|--|
| Mean<br>versus<br>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  $I$ ,  $Y = h(X)$ ;

if  $m_X$  is  $\textcircled{a}$  median of  $X$  (ie,  $\textcircled{205}$ )

if  $m_X = F_X^{-1}(\frac{1}{2})$ , then  $h(m_X)$  is

$\textcircled{a}$  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  $h(x) = ax + b$

$X$  rv with  
mean  $\mu_X$ , SD  $\sigma_X$

Prediction  
~~Feedback~~  
Feedback

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  $\hat{x}$  <sup>206</sup> <sub>x-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.

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

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

a prediction for  $X$ . Small ~~low~~ 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

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

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 (sketch)  $\mathbb{X}$  rv with finite mean  $\mu_{\mathbb{X}}$ ; (209)

let  $m_{\mathbb{X}}$  be (a/the) median of  $\mathbb{X}$ ;

$\rightarrow$  the  $x$  that minimizes  $MAD(x)$  is (a/the) median  $m_{\mathbb{X}}$ . Reminder: why a/the?

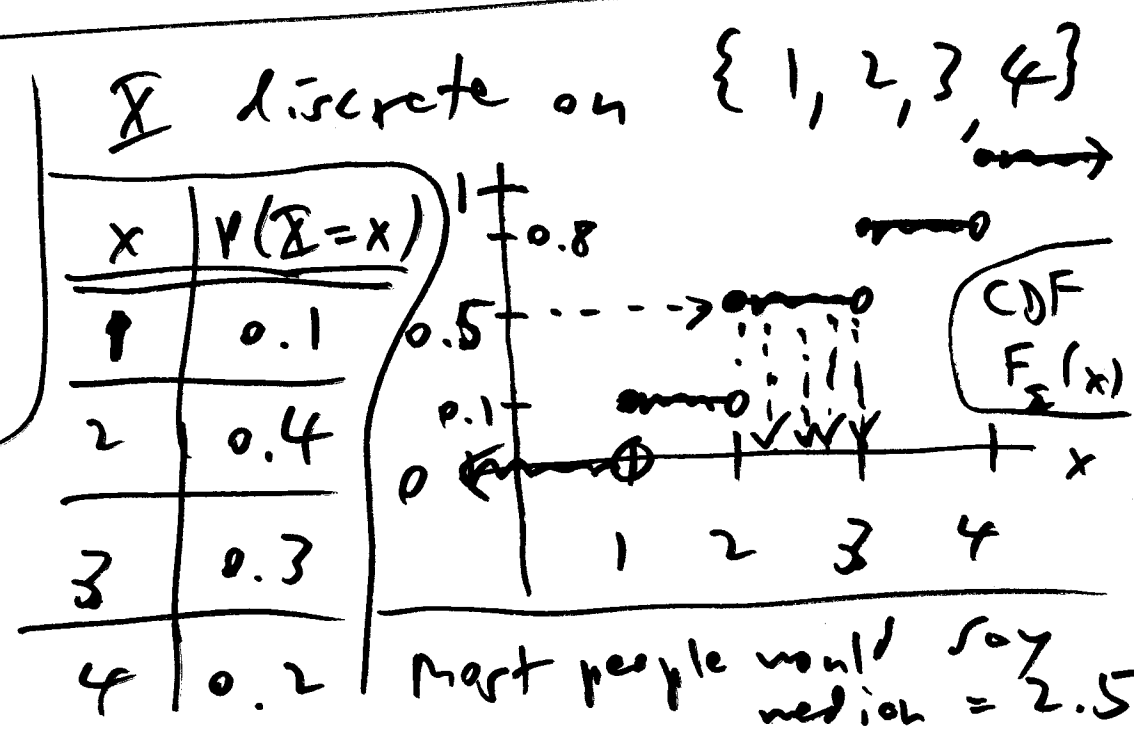
Careful definition of median

$\mathbb{X}$  rv  $\rightarrow$  every number  $m$  such that  $P(\mathbb{X} \leq m) \geq \frac{1}{2}$  and  $P(\mathbb{X} \geq m) \geq \frac{1}{2}$

is a median of the dist. of  $\mathbb{X}$

Example of nonunique median

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



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.

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

Use  $Cov(X, Y)$

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

as long as this expectation exists

Consequences of this definition

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

$$\begin{aligned} \text{so } C(X, Y) &= E(XY) - \mu_X E(Y) - \mu_Y E(X) \\ &= E(XY) - \mu_X \mu_Y - \mu_Y \mu_X + \mu_X \mu_Y \end{aligned}$$

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

(expectation of product - product of expectations)

much easier formula to compute with

② Sufficient condition for  $C(X, Y)$  to 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:  $X$  = education level (years of schooling completed)

$Y$  = yearly income (\$)

Example:

$X$  = temperature <sup>max daily</sup>

in °C <sup>max daily</sup>

$Y$  = humidity (%) <sup>relative</sup>

$C(X, Y)$  comes out in (years) · (\$) (??)

If you change your

mind & measure temperature in °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  $a, d$  are <sup>fixed</sup> constants (23)

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

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

of  $\uparrow$   $\uparrow$  make the association  
between temperature & relative  
humidity seem larger just by switching  
from °C to °F (???)

Easy fix:

Def The process of converting a var  $X$

to standard units (SU) is achieved with

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

(or by  $a = 1/SD(X) < \infty$ , 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

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

$\downarrow$  rho("cov")

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  $a, c \geq 0$  and  $b, d$ ,

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

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

(19 Aug 19)

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

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

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

(22 May 20)

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 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  $\rho$  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

---

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

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

---

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



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

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

④ Important:

if

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

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

(bedrock data science formula)


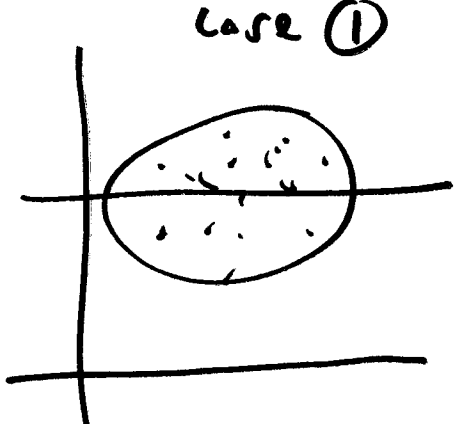
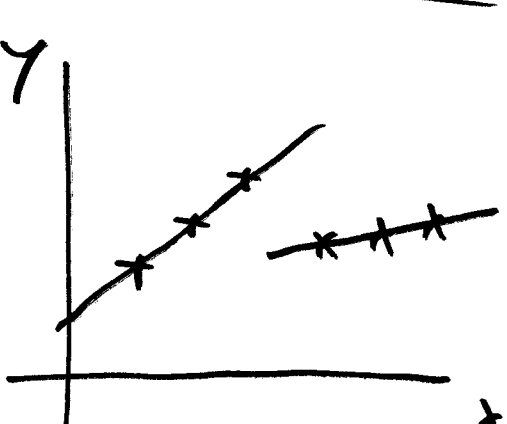
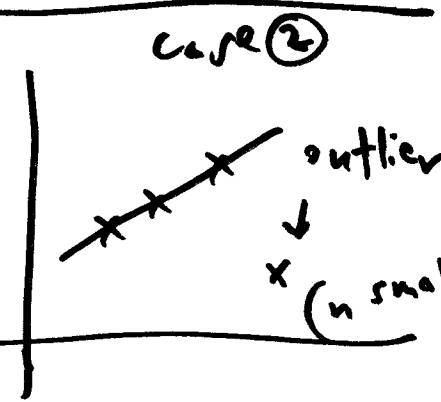
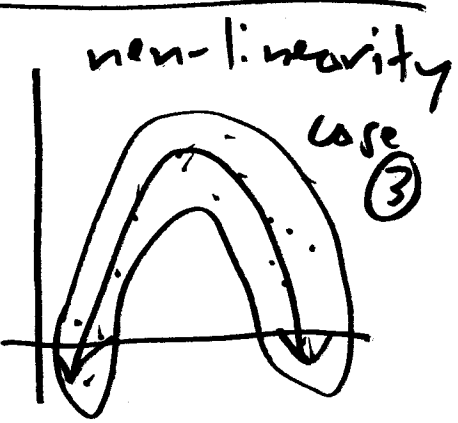
⑤  $\left. \begin{matrix} a, b, c \\ \text{any} \\ \text{constants} \end{matrix} \right\} C(aX, bY + c) = ab C(X, Y)$

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

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

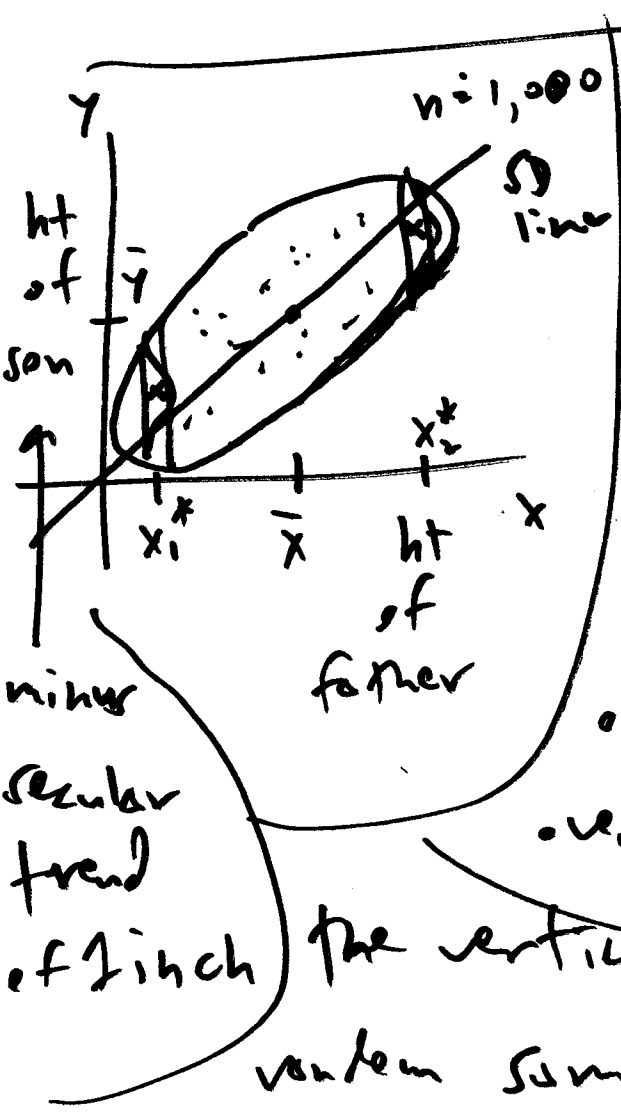
$$V(X-Y) = V(X) + V(Y) - 2C(X, Y)$$

⑥ <sup>(218)</sup>  $\mathbb{R}^n$  such that  $(\mathbb{X}_i, \mathbb{Y}_i)$  uncorrelated  
 for all  $1 \leq i \neq j \leq n \rightarrow$  (then)  $\sqrt{\left(\sum_{i=1}^n \mathbb{R}_i\right)} = \sum_{i=1}^n \sqrt{\mathbb{R}_i}$

| ⑦ $\rho(\mathbb{X}, \mathbb{Y}) = -1$  | $\rho(\mathbb{X}, \mathbb{Y}) = 0$   | $\rho(\mathbb{X}, \mathbb{Y}) = +1$  |
|--|--|--|
|   | <p>Case ①</p>    |    |
| <p>points in scatterplot sample from <math>f_{\mathbb{X}, \mathbb{Y}}(x, y)</math> all fall on line with negative slope (not necessarily -1)</p> | <p>Case ②</p>  <p>Case ③</p> <p>non-linearity</p>  | <p>points in scatterplot sample from <math>f_{\mathbb{X}, \mathbb{Y}}(x, y)</math> all fall on line with positive slope (not necessarily +1)</p> |

(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 ~~plot~~:

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

lect on p. (207): the number  $\hat{w}$  that minimizes the mean squared error,  $E[(\hat{w} - W)^2]$  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.  $X, Y$  r.v.,  $Y$  finite mean  $\rightarrow$

conditional expectation (mean) of  $Y$  given  $X=x$  } =  $E(Y|x)$  is just

the expectation of the conditional distribution

$f_{Y|X}(y|x)$  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  $Y$

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  $X$ . (21/19)

Clinical trial example, continued

$(n_C + n_T)$  people<sup>(A)</sup> 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). (C)

outcome of interest is dichotomous:

let  $\theta$  be the proportion of successes you would have seen if you

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

could have put (everybody in  $P$ ) into your treatment group;  $\theta$  is unknown.

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

Then the rvs  $(S_i | \theta)$  are IID Bernoulli( $\theta$ ) <sup>(23)</sup>

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

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

---

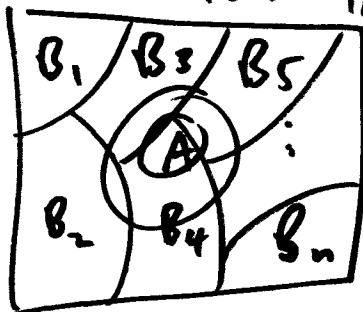
It's meaningful to talk about the conditional expectation rv.  $E(S | \theta) = n_T \theta$  (a linear function of  $\theta$ ),

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  $\rightarrow$

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

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:

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

if ok to interchange order of integration

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



$$= \int_{-\infty}^{\infty} f_X(x) \cdot E(Z|x) dx^*$$

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

Recall that  
continuous  
for any  $v, w$ ,

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

and

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

so  $\textcircled{*}$  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  $\textcircled{1}$  of the  
double expectation theorem; strangely, I  
don't even mention that name, calling it instead  
the LTV for expectations.

I need to postpone examples of these 226  
conditional expectation calculations until  
we've covered more standard distributions.

---

~~Let~~  $X, Y$  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  $V(Y|x) \stackrel{\Delta}{=} E \left\{ [Y - E(Y|x)]^2 \mid x \right\}$   
is called the conditional variance of  $Y = g(X)$   
 $Y$  given  $X = x$ , and the r.v.  $V(Y|X)$  is  
just  $g(X)$ , the conditional variance  
of  $Y$  given  $X$ .

---

The payoff  
from all  
of this

(formalizing Galton's intuition) (227)

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  $X$ .

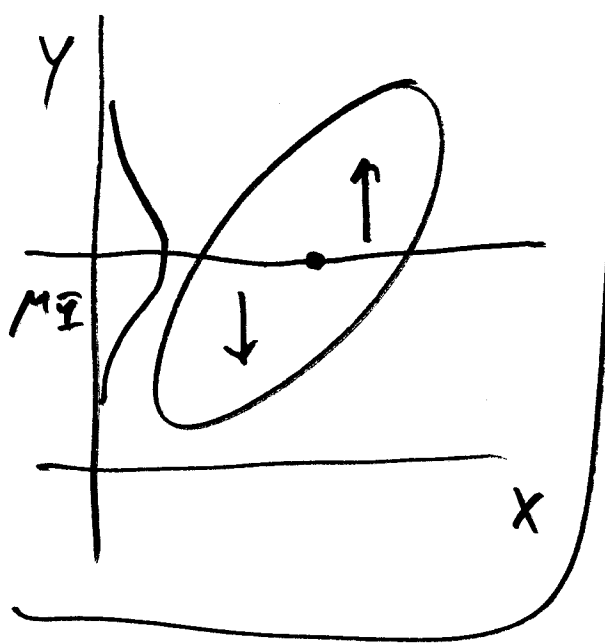
$X, Y$  r.v. such that all of the  
following expressions exist,  $\rightarrow$

$$V(\bar{Y}) = E_X [V(Y|X)]$$

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

(Eve)

Part (2)  
of the  
double  
expectation  
theorem



Imagine a 2-part game!

Stage 1 Predict  $Y$  without knowing  $X$ . Well, if you buty into MSE as your

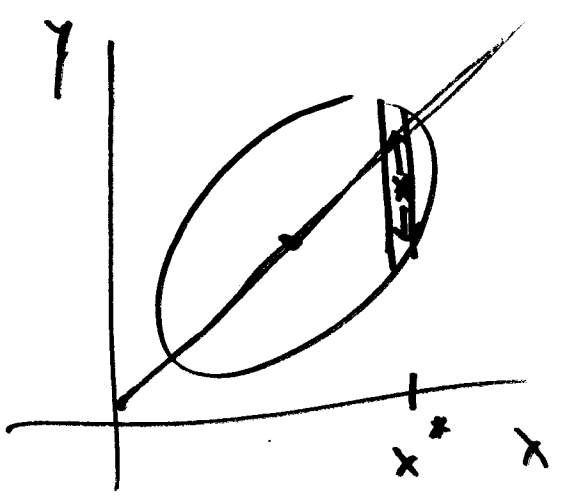
measure of "goodness" of a prediction, we know that you should predict  $\hat{Y}_{no X} = \mu_Y = E(Y)$

and your resulting MSE will be

$$E[(Y - \mu_Y)^2] = V(Y) = \sigma_Y^2$$

Stage 2

observe  $X$ , now predict  $Y$



let's say  $X = x^*$

Then we

know the MSE-optimal

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

and your resulting MSE will be

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

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

namely  $E_X [V(\bar{Y} | X)]$ , 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

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

But since variances are always non-negative,

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

$$E_{\mathcal{X}} [ V(\mathcal{Y} | \mathcal{X}) ] + V_{\mathcal{X}} [ E(\mathcal{Y} | \mathcal{X}) ] \geq E_{\mathcal{X}} [ V(\mathcal{Y} | \mathcal{X}) ]$$

$$V(\mathcal{Y}) \geq \text{MSE of } \hat{\mathcal{Y}}_{\text{no } \mathcal{X}}$$

"E(MSE)"  
of  $\hat{\mathcal{Y}}_{\mathcal{X}}$

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

Another complete switch in subject

Utility

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

A: There is a theory of optimal action under uncertainty; it's called Bayesian decision theory - a concept called utility

is central to this theory. The theory takes its simplest form when comparing ~~gambles~~ gambles

Example  $X$  has discrete PF  $f_X(x) = \begin{cases} \frac{1}{2} & x = -\$350 \\ \frac{1}{2} & x = +\$500 \\ 0 & \text{else} \end{cases}$

Suppose  $X =$  your net gain from gamble (A),

$Y$  has discrete PF  $f_Y(y) = \begin{cases} \frac{1}{3} & y = \$40 \\ \frac{1}{3} & y = \$50 \\ \frac{1}{3} & y = \$60 \\ 0 & \text{else} \end{cases}$

and  $Y =$  your net gain from gamble (B).

Turns out that So is (A) automatically better than (B)?  
 $E(X) = \$75, E(Y) = \$50$

Note that with (B) you're guaranteed to (L32)  
win at least \$40, while (A) has no  
such guarantee; is (A) still automatically  
better for you than (B)? A risk-averse

person would grab (B) quickly; a  
risk-seeking person would probably pick (A).

Evidently something more than just  
computing  $E(X)$ ,  $E(Z)$  is going on.

Def. of utility function

Your utility function  $U(x)$   
is that function which assigns  
to each possible net gain

$-\infty < x < \infty$  a real #  $U(x)$  representing the  
value to you of gaining  $x$ .



Q: If  $x$  is money, why not just use  $u(x) = x$ ? (233)

$u(x) = x$ ?  
(linearity in money)

(A) Well, subtle answer first  
supplied by Daniel Bernoulli (1700 - 1782),  
↙ (Swiss mathematician)  
related to Jacob Bernoulli (1654 - 1705), for  
whom the Bernoulli distribution was named.

Daniel B: If your entire net worth is (say)  
\$10, then the value to you of a new \$1

is much greater than if your entire net  
worth is (say) \$1,000,000; thus the  
utility of money is sublinear (meaning

that it doesn't grow with  $x$  as fast as

$f(x) = x$  does)

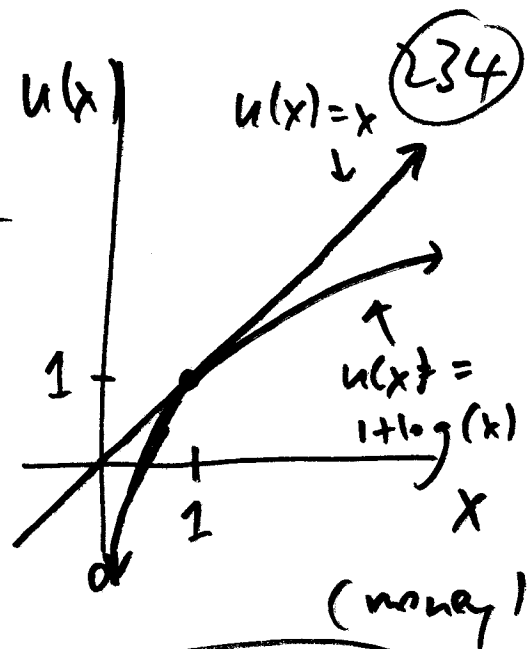
Daniel B proposed one

particular sublinear function for utility.

namely  $u(x) = 1 + \log(x)$   
(for  $x > 0$ )

(Daniel B also invented the  
word utility) (Although

the idea goes back at least  
to Aristotle (384-322 BCE))



Definition

(Principle of  
Expected  
Utility  
Maximization)

You are said to choose  
between gambles by maximizing expected utility (MEU)

if, with  $u(x)$  your utility function,

① you prefer gamble  $\mathbb{X}$  to gamble  $\mathbb{Y}$

if  $E[u(\mathbb{X})] > E[u(\mathbb{Y})]$  and ② you're

indifferent between  $\mathbb{X}$  and  $\mathbb{Y}$  if  $E[u(\mathbb{X})] = E[u(\mathbb{Y})]$

MEU first explored in depth by British (235)

{ mathematician  
philosopher  
economist } Frank Ramsey (1903 - 1930)  
who died at <sup>age</sup> 26 of liver failure.  
(hepatitis)

Theorem / (von Neumann - Morgenstern  
(1947))

John von Neumann  
(1903 - 1957)

Under 4 reasonable axioms,  
MEU is the best you can do.

Hungarian - American

Simple example / Suppose  
You bought

{ mathematician  
physicist  
computer scientist }  
:

a single \$2 ticket in

the Power Ball lottery examined

died at 53 of  
cancer

in ~~the~~ <sup>Take-Home Test</sup> problem 2:  
the drawing on 30 Jul 2016

for which the Grand Prize

Oskar Morgenstern  
(1902 - 1977)

was \$487 million. Let  $X$

German economist  
American

be the amount you will win  
<sup>unknown</sup>

(Thinking about  $X$  before the drawing).

| Match  | $x$                    | $P(X=x)$                                   | $x \cdot P(X=x)$ (236) |
|--------|------------------------|--|------------------------|
| 5w, 1R | \$487,000,000          | $\frac{1}{292,201,338}$                    | \$1.667                |
| 5w, 0R | \$1,000,000            | $\frac{1}{11,688,05352}$                   | 0.086                  |
| 4w, 1R | \$50,000               | $\frac{1}{913,129,18}$                     | 0.055                  |
| 4w, 0R | \$100                  | $\frac{1}{36,525,17}$<br><del>0.0000</del> | 0.003                  |
| 3w, 1R | <del>\$100</del> \$100 | $\frac{1}{14,494,11}$<br><del>0.0000</del> | 0.007                  |
| 3w, 0R | \$7                    | $\frac{1}{579,76}$<br><del>0.0000</del>    | 0.012                  |
| 2w, 1R | \$7                    | $\frac{1}{701,33}$                         | 0.010                  |
| 1w, 1R | 84                     | $\frac{1}{91,98}$                          | 0.043                  |
| 0w, 1R | \$4                    | $\frac{1}{38,32}$                          | 0.104                  |
|        |                        |  | \$1.99 (!)             |

$X$  has 9 possible values  $x$  (discrete),

So  $E(X) = \sum_{\substack{\text{all} \\ 9 \text{ possibilities}}} x \cdot P(X=x) = \$1.99$

**Q:** Before the drawing, someone offers you  $\$x_0$  for your ticket; should you

sell?

**A:** with  $U(x)$  as your utility function, your expected gain if you keep the ticket is  $E[U(X)]$ ; if for you  $U(x) = x$  (utility  $\hat{=}$  money) then

$$E[U(X)] = \$1.99$$

Action 1 (sell): you gain  $\$x_0$  for sure

Action 2 (keep):

your expected utility is  $E[U(X)]$

Under MEU you should sell if  $U(x_0) > E[U(X)]$

If  $U(x) = x$  for you then your optimal action is (sell if offered more than  $\$1.99$ ).

Related but  
different  
problem

On <sup>the</sup> 13 Jan 2016 drawing the 238  
Powerball jackpot was \$1.6 billion.

$X$  = year winnings

$X$  uncertain before  
the drawing

redo calculation on p. 236:  $E(X)$  is  
now \$5.80 or a \$2 ticket

$$\begin{array}{r} \text{new 1st} \\ \text{row in} \\ \text{table is} \\ \hline 1,600,000,000 \\ \hline 292,201,338 \\ \hline = \$5.476 \end{array}$$

(Q.: If  $u(x) = x$  for you,  
under MEU

is it rational to sell all \*

your assets & buy as many lottery  
tickets as possible?

A: Yes, but that's

a silly utility function; to be realistic  
you'd have to subtract from  $x$  the

monetary values <sup>(cost)</sup> to you of the disruption (239)  
of your life that would ensue with action  
(23 May 19)

(\*) A catalog of useful distributions

(Dsch.5) Case 1: Discrete Bernoulli

$X \sim \text{Bernoulli}(p)$ ,  $0 < p < 1$ , if

$$f_X(x) = p^x (1-p)^{1-x} \mathbb{I}_{\text{support}(X)}(x)$$

$$= \begin{cases} p & \text{for } x=1 \\ 1-p & \\ 0 & \text{else} \end{cases}$$

$$E(X) = p$$

$$\psi_X(t) = pe^t + (1-p) \text{ for}$$

all  $-\infty < t < \infty$

$$V(X) = p(1-p)$$

$$SD(X) = \sqrt{p(1-p)}$$

Def If the  $X_i$  in  $X_1, X_2, \dots$  are IID Bernoulli ( $p$ ), then  $(X_1, X_2, \dots)$  are called Bernoulli trials with parameter  $p$ ; if the sequence  $(X_1, X_2, \dots)$  is infinite this defines a Bernoulli (stochastic) process.

Binomial  $X \sim \text{Binomial}(n, p)$  (i.e.,  $X$  follows the Binomial distribution with parameters  $n$  (positive integer) and  $0 < p < 1$ )

$\leftrightarrow f_X(x) = \binom{n}{x} p^x (1-p)^{n-x} \mathbb{I}_{\text{support}(X)}(x)$

Support(X)

Consequences  $X_1, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(p)$

$\rightarrow X = \sum_{i=1}^n X_i \sim \text{Binomial}(n, p)$



$X \sim \text{Binomial}(n, p)$   $E(X) = n \cdot p$  /  $V(X) = n \cdot p \cdot (1-p)$  (24)

$\psi_X(t) = [pe^t + (1-p)]^n$  for all  $-a < t < a$

$SD(X) = \sqrt{np(1-p)}$

Case Study

Supreme Court case

Cartaneda v. Partida (1977)

Grand juries in the U.S. judicial system have 18  
catchment areas: everybody ~~is~~ & over  
 living in the judicial district for that grand  
 jury (a few other minor restrictions)

Hidalgo  
 County,  
 Texas  
 extreme  
 southern  
 border  
 of TX  
 with Mexico

eligible pool was 79.1% Mexican-American

2 1/2 yr period at issue in Supreme

Court case: 220 people called to

serve on grand juries, but only

100 of them were Mexican-American

Q: Prima facie case of discrimination?

Before this 2 1/2 yr period, let  $X$  be your prediction of # of Mexican-Americans among the 220 people

(If) no discrimination,

$X \sim \text{Binomial}(220, 0.791)$   
 $(X | T_1) \rightarrow$

$T_1 = \text{theory}$

$E(X | T_1) = (220)(0.791) = 174.0$

= no discrimination

$SD(X | T_1) = \sqrt{np(1-p)} = 6.0$

Q: If you were

expecting 174 give or take ~~6~~, would you be surprised to see 100?

A: You'd be astonished

Frequentist statistical answer

$P(X \leq 100 | T_1) = 8.0 \cdot 10^{-28}$   
 $T_1$  looks ridiculous

Bayesian statistical answer

Need to compute  $P(T_1 | X = 100)$ , not the other way around (later)

Hypergeometric } A finite population has 243

A elements of type 1 and B elements of type 2; total population size  $(A+B)$ .

---

You choose  $n$  elements at random without replacement from this population (ie, you take a simple random sample (SRS) of size  $n$ )

Let  $X =$  (# elements of type 1 in your sample)

Then (as noted in Take-Home Test 1 problem 2)  $X$  follows the

hypergeometric distribution with

parameters  $(A, B, n)$ .

As we saw

in that problem, the  $P.F.$  of  $X$  is

$$f_{\mathcal{X}}(x | A, B, n) = \frac{\binom{A}{x} \binom{B}{n-x}}{\binom{A+B}{n}} \mathbb{I}[\max\{0, n-B\} \leq x \leq \min\{n, A\}]$$

Support( $\mathcal{X}$ ) (24)

for  $(A, B, n)$  non-negative integers with

$$n \leq A+B$$

Consequences

$$\textcircled{1} E(\mathcal{X}) = n \cdot \frac{A}{A+B}$$

$$\textcircled{2} V(\mathcal{X}) = n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right) \left( \frac{A+B-n}{A+B-1} \right)$$

Note that if

your sampling had been with replacement (i.e., you take an IID sample),  $\mathcal{X}$

would have been Binomial with the

same value of  $n$  and  $p = \frac{A}{A+B}$ ; in

that case  $E(\mathcal{X}) = np = n \frac{A}{A+B}$  and

$$V(\mathcal{X}) = np(1-p) = n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right) \quad (\text{compare})$$

If you let  $T = (A+B)$  be the total # 245  
of elements in the population,

| Sampling method        | mean of $\bar{x}$                | variance of $\bar{x}$  |
|------------------------|----------------------------------|--|
| with repl.<br>(IID)    | $n \left( \frac{A}{A+B} \right)$ | $n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right)$                                |
| without repl.<br>(SPS) | $n \left( \frac{A}{A+B} \right)$ | $n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right) \left( \frac{T-n}{T-1} \right)$ |

$0 \leq \alpha = \frac{T-n}{T-1} \leq 1$  is called the finite

population correction  
(21 Aug 19)

3 special cases  
worth considering

(a)  $(n=1) \alpha=1 \leftrightarrow$  SPS = IID with only  
1 element  
sampled

(b)  $(n=T) \alpha=0 \leftrightarrow$  If you  
exhaust the entire population <sup>with</sup> SPS,  
you have no uncertainty left.

(c) ( $n$  fixed,  $T \uparrow$ )  $d \xrightarrow{1} \leftrightarrow$  with a 246  
 small sample from a large population,

$$S_{ij} = IID$$

Poisson ( $\lambda > 0$ )  $X \sim \text{Poisson}(\lambda)$

$\leftrightarrow X$  has  $\overset{M}{PF}$   $f_X(x) = \frac{\lambda^x e^{-\lambda}}{x!} \mathbb{I}_{\{0, 1, \dots\}}(x)$   
support of  $X$

$$E(X) = \lambda$$

$$V(X) = \lambda$$

thus for the Poisson dist.

$$\frac{V(X)}{E(X)} = 1 \quad \text{Def.} \quad \text{If } E(X) \text{ and } V(X)$$

$$\psi_X(t) = e^{\lambda(e^t - 1)}$$

$$-\infty < t < \infty$$

both exist and  $E(X) \neq 0$ ,

$\frac{V(X)}{E(X)}$  is called the

variance-to-mean ratio

(VTMR)

→ because

The Poisson can be unrealistic as a consequence of its VTMR of 1,

many rvs that represent counts of 247  
occurrences of events in time intervals  
of fixed length have  $VMR > 1$ .

---

The Poisson & Binomial distributions  
both count the number of "successes"  
in a process unfolding in time, so  
it should not be surprising to find  
out that these 2 dist. are related:

---

when  $\begin{pmatrix} n \text{ is large} \\ p \text{ is close to } 0 \end{pmatrix}$ ,  $\text{PMF of Binomial}(n, p) \doteq \text{PMF of Poisson}(n \cdot p)$

---

Theorem  $n$  positive integer,  $0 < p < 1$   $X \sim \text{Binomial}(n, p)$

---

$\lambda > 0$ ,  $Z \sim \text{Poisson}(\lambda)$  / Choose any sequence

$\{p_n\}_{n=1}^{\infty}$  of values between 0 and 1 with

$\lim_{n \rightarrow \infty} n \cdot p_n = \lambda$

PMF Binomial

Then  $f_X(x | n, p_n) \rightarrow_{n \rightarrow \infty}$

Poisson process, revisited

Def

$f_X(y | \lambda)$   
PMF Poisson

A Poisson process with rate  $\lambda$  per unit (or space, or volume, or...)

time, is a stochastic process with two

properties:

(a) # arrivals in every interval of time of length  $t \sim \text{Poisson}(\lambda t)$

(b) #s of arrivals in all disjoint (non-overlapping) time intervals are independent

Core Study

~~Parasitic~~  
Parasitic protozoa in drinking water

There's a kind of parasitic



organism called cryptosporidium that's (249)  
capable of getting into the public drinking  
water supplies; at one stage in their life  
cycle they're called ooocysts.

They can make  
people sick at a concentration of only  
1 ooocyst per 5 liters = 1.3 gallons of water

One problem is that it can be hard to detect  
these ooocysts with water filtration.

Suppose  
that, in the water supply of your city,  
ooocysts occur according to a Poisson process  
with rate  $\lambda$  ooocysts per liter, & that  
the filtering system your water utility  
company uses can capture all the ooocysts  
in a water sample but only has

probability  $p$  of detecting each oocyst 250

that's actually there. (Counting events are independent)

Let  $\underline{Y}$  = # oocysts in  $t$  liters of water,  
and  $\underline{X}_i = \begin{cases} 1 & \text{if oocyst } i \text{ gets counted} \\ 0 & \text{else} \end{cases}$

$\underline{X}$  = # counted oocysts | Then  $(\underline{X} | \underline{Y} = y) = \sum_{i=1}^y \underline{X}_i$

under these assumptions,  $(\underline{X} | \underline{Y} = y) \sim \text{Binomial}(y, p)$

Q: what's the dist. of  $\underline{X}$ ? | A | By the

law of total probability

$$f_{\underline{X}}(x) = P(\underline{X} = x) = \sum_{y=0}^{\infty} P(\underline{Y} = y) P(\underline{X} = x | \underline{Y} = y)$$

for all  $x = 0, 1, \dots$

in which  $P(\underline{Y} = y) = \frac{(\lambda t)^y e^{-\lambda t}}{y!}$  for  $y = 0, 1, \dots$