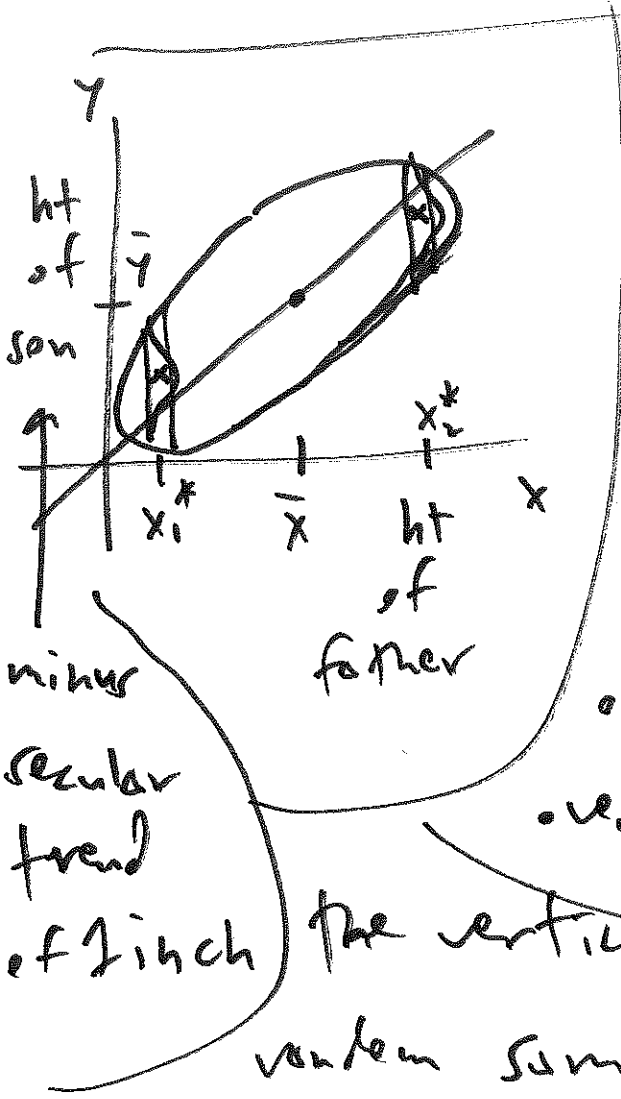


(21 Aug 19)  
 Conditional  
 Expectation

$X, Y$  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  $x_2^*$ .

~~over~~ The points in the strip over  $x_2^*$  are a 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 (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 "close" 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.  $E, E v, Y$  finite mean  $\rightarrow$

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

the expectation of the conditional distribution

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

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

for continuous  $(Y|X=x)$

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

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  $X$ .

Clinical trial example, continued

$(n_C + n_T)$  people<sup>(a)</sup> 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  $\theta$  be the proportion of successes you would have seen if you could have put (everybody in  $P$ ) into your treatment group;  $\theta$  is unknown.

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

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

---

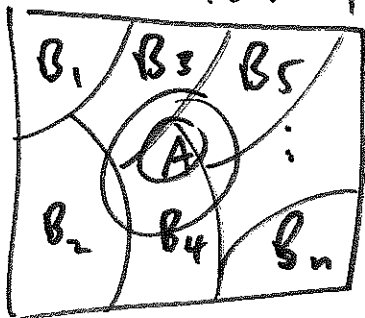
It's meaningful to talk about the conditional expectation rv.  $E(S | \theta) = n_T \theta$  (a linear function of  $\theta$ ),

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  $\rightarrow$

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

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:

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

if ok to interchange order of integration

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

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

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

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

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 226  
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  $V(Y|x) \stackrel{\Delta}{=} E \left\{ \left[ Y - E(Y|x) \right]^2 \mid x \right\}$   
is called the conditional variance of  $Y$  given  $X=x$ .  
 $\overset{=g(x)}{\uparrow}$   
 $\overset{=g(x)}{\uparrow}$   
I give  $X=x$ , 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  $E(Y - \hat{Y})^2 = E\left\{[Y - d(X)]^2\right\}$

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

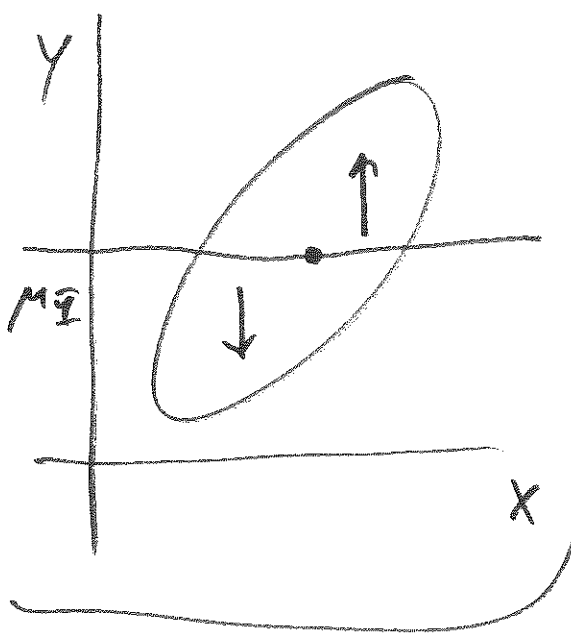
$X, Y$  r.v. such that all of the  
following expressions exist,  $\rightarrow$

$$V(Y) = E_X[V(Y|X)]$$

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

(Eve)

Part (2)  
of the  
double  
expectation  
theorem



Imagine a 2-part game!

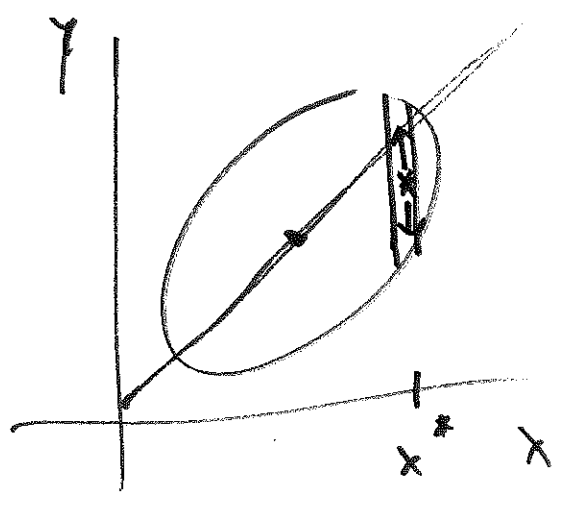
Step 1) Predict  $Y$  without knowing  $X$ . Well, if you but into MSE as your

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

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

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

(229)

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

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

namely  $E_{\mathcal{X}} [V(\mathcal{Y} | \mathcal{X})]$ , 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

$$\underbrace{V(\mathcal{Y})}_{\substack{\uparrow \\ \text{MSE of} \\ \hat{\mathcal{Y}}_{\text{no } \mathcal{X}}}} = E_{\mathcal{X}} \left[ \underbrace{V(\mathcal{Y} | \mathcal{X})}_{\substack{\text{"E(MSE)" of} \\ \uparrow \\ \hat{\mathcal{Y}}_{\mathcal{X}} = E(\mathcal{Y} | \mathcal{X})}} \right] + \underbrace{V_{\mathcal{X}} [E(\mathcal{Y} | \mathcal{X})]}_{\substack{\text{variance of} \\ \text{the conditional} \\ \text{mean}}}$$

But since variances are always non-negative,

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

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

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

"E(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?

A: There is a theory of optimal actions 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)

and  $Y =$  your net gain from gamble (B).  $Y$  has discrete PF  $f_Y(y) = \begin{cases} \frac{1}{3} & y = \$40 \\ \frac{1}{3} & y = \$50 \\ \frac{1}{3} & y = \$60 \\ 0 & \text{else} \end{cases}$

Turns out that So is (A) automatically better than (B)?  
 $E(X) = \$175, E(Y) = \$50$

Note that with (B) you're guaranteed to (232)  
win at least 84%, 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

$-∞ < x < ∞$  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)

$u(x) = x$ ?  
(utility = money)

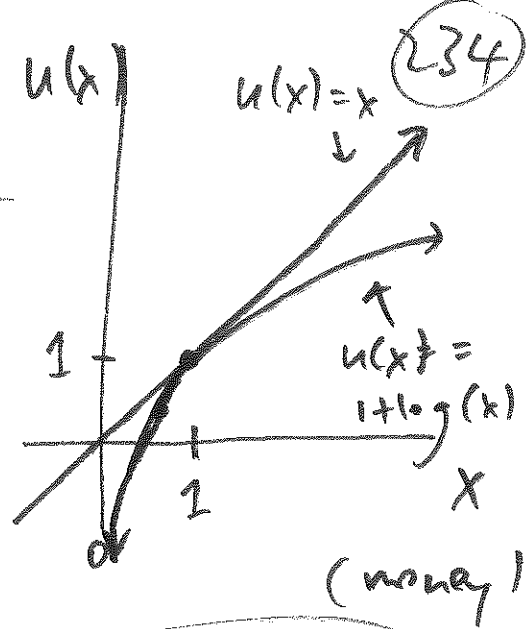
A: lovely, subtle answer first

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

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 Brit. (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  
in ~~Take-Home~~ Test  
in ~~problem~~ problem 2:

died at 53 of  
cancer

the drawing on 30 Jul 2016  
for which the Grand prize  
was \$487 million. Let  $X$   
be the <sup>unknown</sup> amount you will win

Oskar Morgenstern  
(1902 - 1977)

Game-theorist  
American

(think 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	$\frac{1}{11,688,053.52}$	0.086
4w, 1R	\$50,000	$\frac{1}{913,129.18}$	0.055
4w, 0R	\$100	$\frac{1}{36,525.17}$ <del>0.0049</del>	0.003
3w, 1R	<del>\$100</del> \$100	$\frac{1}{14,494.11}$ <del>0.0069</del>	0.007
3w, 0R	\$7	$\frac{1}{579.76}$ <del>0.0017</del>	0.012
2w, 1R	\$7	$\frac{1}{701.33}$	0.010
1w, 1R	\$4	$\frac{1}{91.98}$	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.

$X$  = your winnings

$X$  uncertain before  
the drawing

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

new 1st  
row in  
table is  
 $\frac{1,600,000,000}{292,201,338}$   
 $\approx 5.476$

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

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 value <sup>(cost)</sup> to you of the disruption (239)  
of your life that would ensue with action

(\*) A catalog of useful distributions

(Dsch.5) Case 1: Discrete Bernoulli

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

$$f_X(x) = p^x (1-p)^{1-x} \cdot \underbrace{I_{\{0,1\}}(x)}_{\text{support}(X)}$$

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

$$E(X) = p$$

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

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 240  
 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  $0 < p < 1$ )

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

$\text{Support}(X) = \{0, 1, \dots, n\}$

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

Case Study  
~~Cartaneda~~ Cartaneda v. Partida (1977)

Grand juries in the U.S. judicial system have  
 catchment areas: everybody <sup>18</sup> & 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,

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

$T_1 = \text{theory}$

$E(X | T_1) = \binom{n}{p} = (220)(0.791) = 174.0$

= no discrimination

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

Q: If you were

expecting 174 give or take ~~6~~ 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  
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 (i.e.,  
you take a simple random sample (SRS)  
of size n)

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

~~Then~~ (as noted in  
type-Hone  $T_2^+$   
~~problem 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

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

support( $\mathbb{X}$ ) (244)

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

$$n \leq A+B$$

Consequences

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

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

Note that if

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

would have been Binomial with the

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

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

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

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

Sampling method	mean	variance
with repl. (IID)	$n \left( \frac{A}{A+B} \right)$	$n \left( \frac{A}{A+B} \right) \left( \frac{B}{A+B} \right)$
without repl. (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

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 with SPS, you have no uncertainty left.

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

$SD = \sqrt{E}$

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

$\leftrightarrow X$  has <sup>M</sup>PF  $f_X(x) = \frac{\lambda^x e^{-\lambda}}{x!} \frac{I_{\{0, 1, \dots\}}(x)}{\text{support of } X}$

$E(X) = \lambda$

$V(X) = \lambda$

Thus for the Poisson dist.

$\frac{V(X)}{E(X)} = 1$  Def. If  $E(X)$  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 <sup>can</sup> 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  $VTNR > 1$ .

---

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

---

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

---

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

$\{p_n\}_{n=1}^{\infty}$  of values between 0 and 1 with (248)

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

Then  $f_X(x | n, p_n) \rightarrow$

Poisson process,  
revisited

Def

$$f_X(y | \lambda)$$

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~~  
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 <sup>(250)</sup>

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

Let  $\underline{Y}$  = <sup>actual</sup> # 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  $y = 0, 1, \dots$



and  $P(X=x | Y=y) = \binom{y}{x} p^x (1-p)^{y-x}$  (251)

Notice that if  $X=x$ ,  $Y \geq x$  because the <sup>actual</sup> number of oocysts ( $Y$ ) has to be at least as large as the number of oocysts detected ( $X$ ).

$$f_X(x) = \sum_{y=x}^{\infty} \binom{y}{x} p^x (1-p)^{y-x} \frac{(\lambda t)^y e^{-\lambda t}}{y!}$$

After a careful calculation you get;

$$= \frac{e^{-p\lambda t} (p\lambda t)^x}{x!}$$

i.e.,

$X \sim \text{Poisson}(p\lambda t)$ :  
 losing a proportion  $(1-p)$  of the oocysts to faulty counting just lowers the rate of the Poisson process from  $\lambda$ /liter to  $\lambda \cdot p$ /liter (makes excellent sense).

In practice oocysts are hard to detect <sup>252</sup><sub>t</sub>:

$p$  is small (not far from 0). Q: How

much water <sup>(t liters)</sup> do you need to filter to achieve  $P(\text{at least 1 oocyst detected}) \geq 1 - \alpha$

for small  $\alpha$ ? A: Not hard to work out

$$P(\text{at least 1 detected}) = 1 - P(\text{none detected})$$

$$= 1 - P(X=0) = 1 - e^{-p\lambda t} \geq 1 - \alpha$$

$$\Leftrightarrow \alpha \geq e^{-p\lambda t} \Leftrightarrow \ln \alpha \geq -p\lambda t \Leftrightarrow$$

$$t \geq \frac{-\ln \alpha}{p\lambda}$$

Example)  $\alpha = .01$ ,  $p = 0.1$ ,  
 $\lambda = 0.2 / \text{liter}$  (1 per 5 liters)

to achieve  $p \sim 99\%$ ,

$t$  has to be at least

230.3 liters.

↓  
minimum  
sickness  
level

# Negative Binomial Distribution

You're watching a potential <sup>253</sup> endless sequence of Bernoulli trials with constant success

probability  $p$ .

Let  $X$  = # failures before  $r$ <sup>th</sup> success

You can show that  $X$  follows the Negative Binomial dist:

its PF is  $f(x | r, p) = \binom{r+x-1}{x} p^r (1-p)^x$

with parameters  $(r, p)$

The name comes from the fact that, when you watch a sequence of Bernoulli trials with constant <sup>unknown</sup> success probability  $p$  unfolded, there are two different ways to

estimate  $p$ : decide ahead of time to (254)  
(known constant)  
sample  $n$  success/failure trials, and  
record the (random) #  $S$  of successes  
you see (from which a reasonable  
estimate would be  $\hat{p}_B = \frac{S}{n}$  ← Binomial)

---

(or) decide ahead of time that you're  
going to sample until you've seen  $s$   
(known constant) successes & record the  
(random) # of trials  $N$  needed  
to accumulate that many successes  
(from which a reasonable estimate  
would be  $\hat{p}_{NB} = \frac{s}{N}$  ← Negative Binomial).

Special  
Case of  
Negative  
Binomial

Set  $r=1$  and record the  $(255)$   
number  $X$  of failures until  
the first success:  $X$  is  
said to follow the

Geometric ( $p$ ) distribution, with

$$P\{X = x\} = p(1-p)^x \mathbb{1}_{\{0, 1, \dots\}}(x)$$

(parameter  $p$ )

~~Con~~ source  $X_1, \dots, X_n$  IID Geometric( $p$ )

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

This is a direct analogue to the

Bernoulli/Binomial story:  $X_1, \dots, X_n$  IID

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

$X \sim \text{Negative Binomial}(r, p)$

256

$$\psi_X(t) = \left[ \frac{p}{1 - (1-p)e^t} \right]^r \text{ for } t < \log\left(\frac{1}{1-p}\right)$$

from which  $E(X) = \frac{r(1-p)}{p}$ ,  $V(X) = \frac{r(1-p)}{p^2}$

Consequence

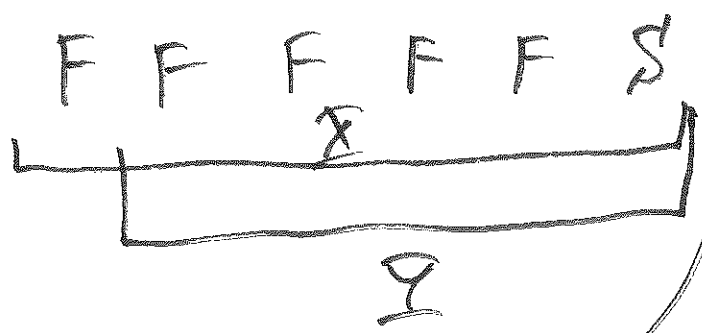
$X \sim \text{Geometric}(p) \rightarrow$

$\begin{cases} k \\ t \end{cases}$  both non-negative integers

$$P(X = k+t \mid X \geq k) = P(X = t)$$

this is called the memoryless property of the Geometric distribution, and it turns out that this is the only

discrete distribution with this property. (257)



$X = \#$  failures until first success = 5 (here)

$Y = \#$  failures, starting at trial  $(k+1)$  until next success ( $= 2$  here,  $\sim 4$  here)

Then  $Y$  has

the same dist. as  $X$  and is independent of what happened on the first  $k$  trials, i.e., "the process has no memory".

Case 2: Important Continuous Distributions

Normal (Gaussian) Distribution

$X \sim \text{Normal}(\mu, \sigma^2)$  mean  $\mu$  variance  $\sigma^2 < \infty$

PDF  $\int$

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

The Normal dist. is the single most important dist. in all of probability & statistics, mainly for 2 reasons:

① many observable random processes have dist. shapes that are close to the "bell curve" (Normal PDF), and

② the Central Limit Theorem (CLT), which we'll examine soon.

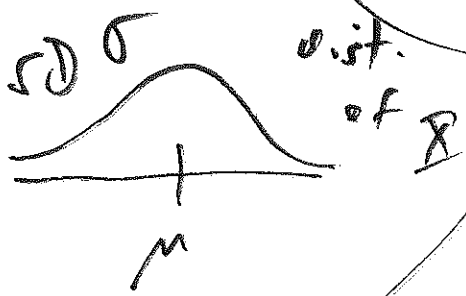
Properties of the Normal Dist.

$$X \sim \text{Normal}(\mu, \sigma^2)$$

$$E(X) = \mu$$

$$V(X) = \sigma^2, \quad SD(X) = \sigma$$

$$\psi_X(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)$$



dist. of X

(center of symmetry)  
mean  
median  
mode  
=  $\mu$



Consequences ①  $X \sim \text{Normal}(\mu, \sigma^2)$ , (259)

$Y = aX + b$ , ( $a \neq 0$ ) fixed constants  $\rightarrow$

$Y \sim \text{Normal}(a\mu + b, a^2\sigma^2)$ .

In other words, Normality is preserved under linear transformations Def.

The Normal dist. with mean  $\mu = 0$  and SD  $\sigma = 1$  is called the standard normal dist.

The PDF of  $X \sim \text{Normal}(0, 1)$  is

$\phi_X(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$  and its  
 $\rightarrow$  phi (lower-case)

CDF is  $\Phi(x) = \int_{-\infty}^x \phi_X(t) dt$   
 $\rightarrow$  uppercase phi  $\Phi$  23  
Aug 17