

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  $V(Y|x) \stackrel{\Delta}{=} E \left\{ [Y - E(Y|x)]^2 \mid x \right\}$   
is called the conditional variance of  $Y$  given  $X = x$ .

$Y$  given  $X = x$ , and the r.v.  $V(Y|X)$  is  
just  $V(Y|X)$ , 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  $E(Y - \hat{Y})^2 = E\{[Y - d(X)]^2\}$

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,  $\rightarrow$

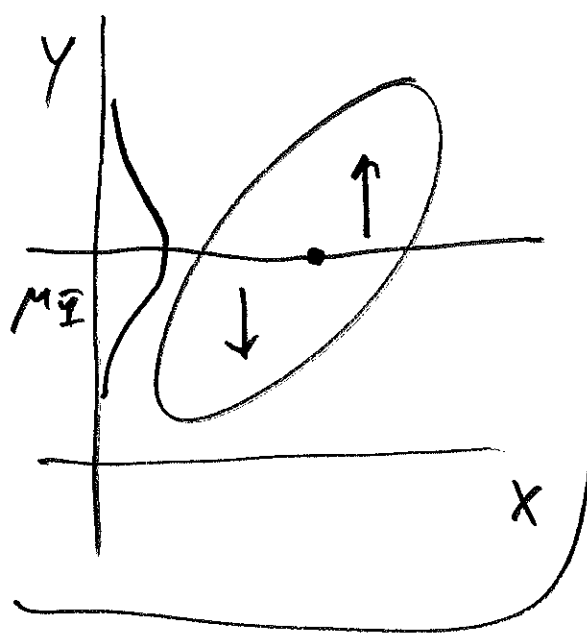
$$V(Y) = E_X [V(Y|X)]$$

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

(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  $\hat{Y}_{no X} = \mu_Y = E(Y)$

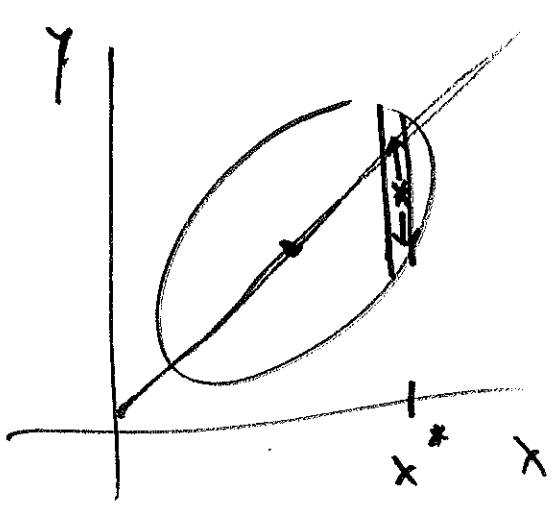
and your resulting MSE will be

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

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

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

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

namely  $E_X [V(Y | X)]$ , 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,

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

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

$$V(Y) \geq E(\text{MSE})$$

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

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

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

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

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

Note that with (B) you're guaranteed to 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

$-a < x < a$  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),  
(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)

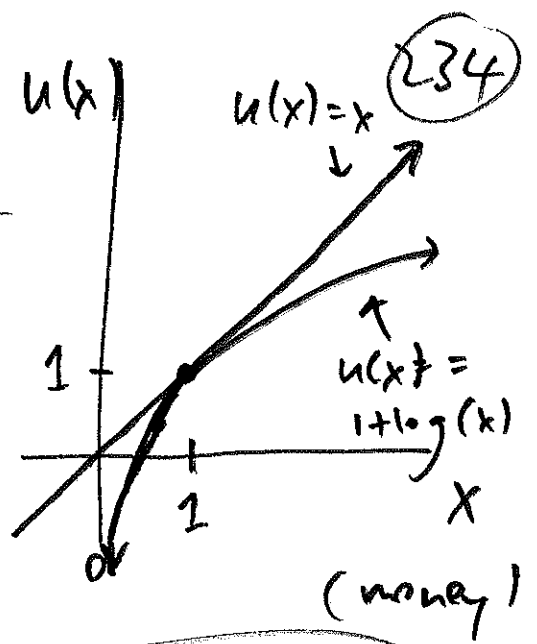
Daniel B proposed one particular sublinear function for utility,



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  
{ mathematician  
physicist  
computer scientist }  
:  
died at 53 of  
cancer

Simple example) Suppose  
you bought

a single \$2 ticket in  
the power ball lottery examined  
in ~~Take-Home~~ ~~Test~~ problem 2:

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

(think about  $X$  before the drawing).

Oskar Morgenstern  
(1902 - 1977)  
German economist  
American

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>0.0027</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, 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

necessary 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

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

$\mathbb{I}_{\{0,1\}}(x)$

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

$$E(X) = p$$

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

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

all  $-\infty < t < \infty$

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

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

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

$X \sim \text{Binomial}(n, p)$   $E(X) = n \cdot p$  /  $V(X) = n \cdot p \cdot (1-p)$  (241)

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

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

Case Study Supreme Court case  
~~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,

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

$E(X | T_1) = \binom{n}{k} p^n = (220)(0.791) = 174.0$  = 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  
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  
type-Homework  
~~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 ①  $E(\mathbb{X}) = n \cdot \frac{A}{A+B}$

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,  $T \uparrow$ )  $\leftrightarrow$  with a small sample from a large population,

SRS = IID

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

$\leftrightarrow X$  has 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.

$\frac{V(X)}{E(X)} = 1$  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 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  $VPMR > 1$ .

---

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

$$f_Y(y | \lambda)$$

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

Case 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  $\mathcal{Y} =$  <sup>actual</sup> # oocysts in  $t$  liters of water,  
and  $\mathcal{X}_i = \begin{cases} 1 & \text{if oocyst } i \text{ gets counted} \\ 0 & \text{else} \end{cases}$

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

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

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

law of total probability

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

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

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