

DS (and some other people) call Method 2 (177) the law of the Unconscious Statistician,

because Method 2 looks like a definition but it actually is a theorem (difficult) (in full generality) (measure theory: pushforward measure, ...)

Example  $X \sim \text{Exponential}(\lambda)$  ( $\lambda > 0$ )  
 $E(X) = \frac{1}{\lambda}$  (integrate by parts twice)

$Y = X^2$   
 $E(Y) = \int_0^{\infty} x^2 \lambda e^{-\lambda x} = \frac{2}{\lambda^2}$

Notice that

$$E(X^2) \neq [E(X)]^2$$
$$\frac{2}{\lambda^2} \neq \left(\frac{1}{\lambda}\right)^2$$

The only functions  $Y = h(X)$  for which  $E[h(X)] = h[E(X)]$  are linear:  $h(x) = a + bx$ , as we'll see later

~~(scribble)~~

Properties of  $E(\underline{Y})$

① If  $\underline{Y} = a \underline{X} + b$  then

$$E(\underline{Y}) = a E(\underline{X}) + b \quad \left( \begin{array}{l} \text{assuming} \\ E(\underline{X}) \\ \text{exists} \end{array} \right)$$

② If you can find a constant  $a$  with  $P(\underline{X} \geq a) = 1$  then (naturally enough)  $E(\underline{X}) \geq a$ ; if  $b$  exists with  $P(\underline{X} \leq b) = 1$  then  $E(\underline{X}) \leq b$ .

③ If  $\underline{X}_1, \dots, \underline{X}_n$  are  $n$  rvs, each with finite  $E(\underline{X}_i)$ , then  $E\left(\sum_{i=1}^n \underline{X}_i\right) = \sum_{i=1}^n E(\underline{X}_i)$ ,

④ and  $E\left[\sum_{i=1}^n (a_i \underline{X}_i + b)\right] = \sum_{i=1}^n a_i E(\underline{X}_i) + b$ .

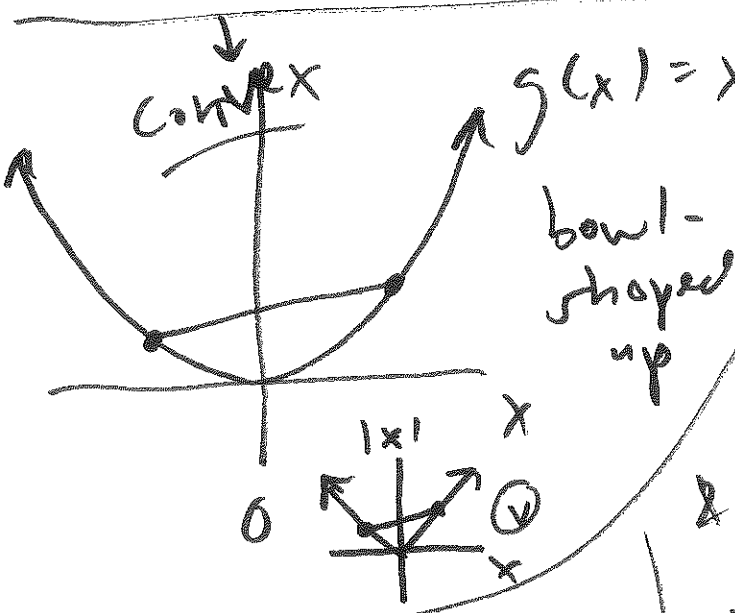
for all constants  $(a_1, \dots, a_n)$  and  $b$ .

⑤ Def. A function  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  (this

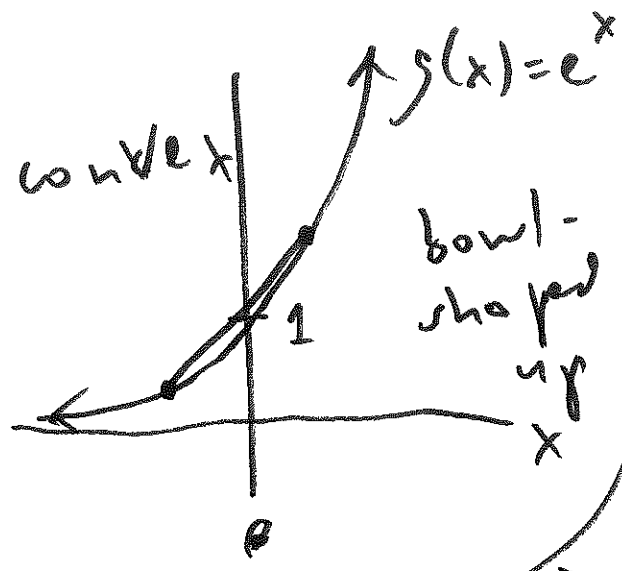
means that  $g(\underline{x}) = z$  is convex  
 $\begin{array}{c} \text{real \#s} \\ \swarrow \quad \searrow \\ \underline{x} = (x_1, \dots, x_n) \end{array}$

if for every  $0 < \alpha < 1$  and every

$$\underline{x} \text{ and } \underline{y}, \quad \underline{g[\alpha \underline{x} + (1-\alpha)\underline{y}]} \leq \alpha \underline{g(\underline{x})} + (1-\alpha)\underline{g(\underline{y})}$$

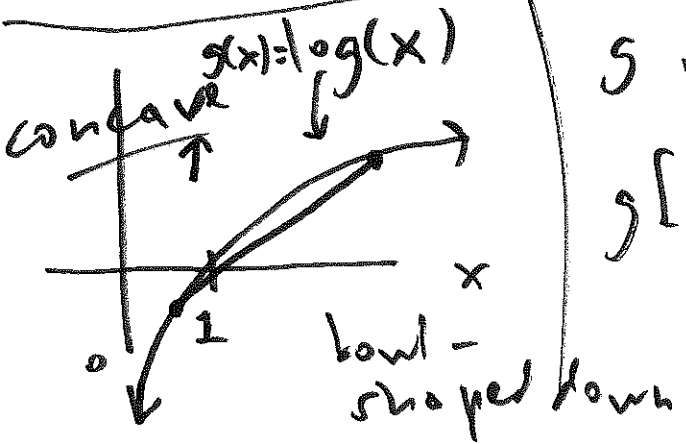


Graphical version of this: pick any two points on the function & connect them with a line segment; the function is convex if the line segment lies <sup>entirely</sup> above the function except at the endpoints.



$g$  is concave if

$$\underline{g[\alpha \underline{x} + (1-\alpha)\underline{y}]} \geq \alpha \underline{g(\underline{x})} + (1-\alpha)\underline{g(\underline{y})}$$



Def. The expectation of a random vector

$\underline{X} = (X_1, \dots, X_n)$  is  $E(\underline{X}) \triangleq [E(X_1), \dots, E(X_n)]$

(a)  $g$  convex,  $\underline{X}$  random vector with finite  $E(\underline{X}) \rightarrow E[g(\underline{X})] \geq g[E(\underline{X})]$ . Jensen's Inequality

(b)  $g$  concave  $\rightarrow E[g(\underline{X})] \leq g[E(\underline{X})]$ .

(attributed to Johan Jensen, (1859-1925), Danish mathematician & engineer)

Applications of (3)

Suppose that  $X_1, \dots, X_n \overset{IID}{\sim} \text{Bernoulli}(p)$ .

Then  $E(X_i) = 0 \cdot \underset{P(X=0)}{\uparrow} (1-p) + 1 \cdot \underset{P(X=1)}{\uparrow} p = p$  and

$E(\sum_{i=1}^n X_i) = \sum_{i=1}^n E(X_i) = np = \text{mean of Binomial}(n, p)$

Expectation  
of a product  
when the  
 $X_i$  are  
independent

$X_1, \dots, X_n$  independent r.v. each with  
finite  $E(X_i) \rightarrow$  (181)

$$E\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n E(X_i)$$

Contrast this with a sum:  $E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i)$   
whether the  $X_i$  are independent or not;

$$E\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n E(X_i) \text{ only when the } X_i$$

are independent.

Example

You have

a (Brita) water filter that you use to  
improve the taste of Santa Cruz water.

How much better would the filter do

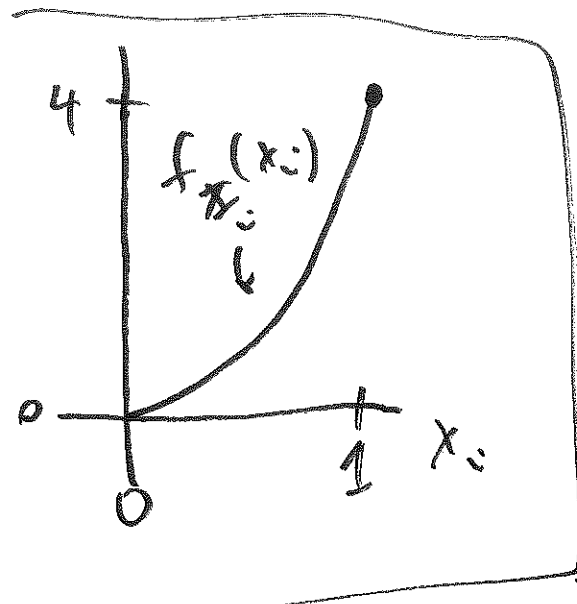
its job if you filtered the water twice  
instead of once?

$X_1$  = proportion of bad stuff removed in the 1<sup>st</sup> filtering (18%)

$X_2$  = proportion removed in 2<sup>nd</sup> filtering of what was left from 1<sup>st</sup> filtering

Reasonable to assume that  $X_1, X_2$  are independent; suppose they're IID with

common PDF  $f_{X_i}(x_i) = \begin{cases} 4x^3 & 0 < x < 1 \\ 0 & \text{else} \end{cases}$



(sensible shape)

Let  $Y$  = proportion of original bad stuff remaining after 2 filtrations =  $(1-X_1)(1-X_2)$

Then  $E(Y) = E[(1-X_1)(1-X_2)] \stackrel{\text{independence}}{=} E[(1-X_1)] \cdot E[(1-X_2)]$

$X_1, X_2$  independent

$\Leftrightarrow (1-X_1), (1-X_2)$  independent too

$E(1-X_1) \stackrel{\text{identical distribution}}{=} E(1-X_2) \triangleq \mu$ ;

then  $E(Y) = \mu^2$ .

$$\mu = E(1 - X_i) = \int_0^1 (1 - x_i) 4x_i^3 dx_i = 0.2, \quad (183)$$

so 80% of bad stuff expected to be removed in 1<sup>st</sup> filtering;  $E(Q) = \mu^2 = 0.04$ , so expect only 4% of bad stuff to remain after 2 filterings.

(6) Suppose  
(9)

$X$  is a discrete rv with possible values  $0, 1, 2, \dots$ ; then  $E(X) = \sum_{n=0}^{\infty} P(X \geq n)$ .  
 $n=0 \leftarrow ?$

(b) If  $X$  is a continuous rv with possible values  $(0, \infty)$ , then  $E(X) = \int_0^{\infty} [1 - F_X(x)] dx$ ,  
and CDF  $F_X(x)$ ,

Example of b(9)

I throw a dart at a dartboard repeatedly, trying to get a bullseye (success).

$X = \#$  of throw on which I first succeed.

(Ex. throws FFS  $\rightarrow X=3$ ) Suppose that my

F = failure  
S = success

success probability is constant across the throws and equals  $p$ , & throws are independent.

Then  $E(X)$  should be inversely related to  $p$ :

The worse I am, the longer I expect the process to take;  $E(X) = ?$

At least 1 throw

always required so  $P(X \geq 1) = 1$ ; for  $n > 1$

(at least  $n$  throws required)  $\leftrightarrow$  (none of the first  $(n-1)$  throws succeeded)

so  $P(X \geq n) = (1-p)^{n-1}$  and arithmetic series

$$E(X) = \sum_{n=1}^{\infty} (1-p)^{n-1} = 1 + (1-p) + (1-p)^2 + \dots$$

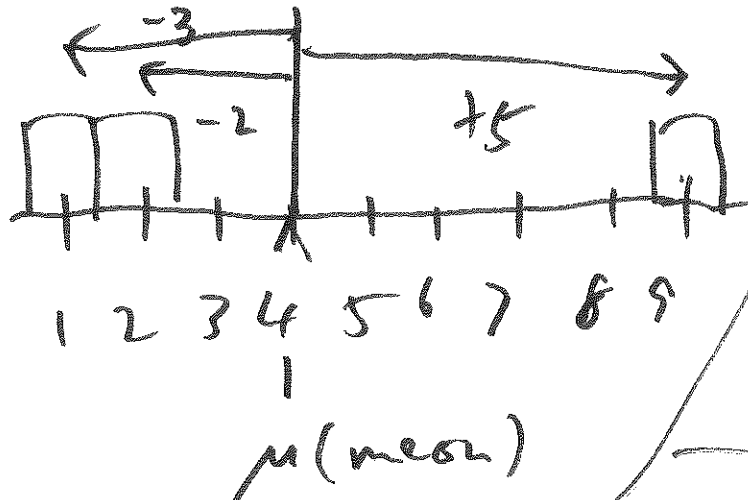
$$= \frac{1}{1-(1-p)} = \frac{1}{p}$$

(inverse relation  $\checkmark$ )

If I'm terrible (e.g.  $p = .01$ ) I expect to succeed on the  $\frac{1}{.01} = 100^{\text{th}}$  throw.



Variance  
and  
standard  
deviation



185

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

mean  $\bar{x} = \mu$

$X$  discrete rv, uniform  $\{1, 2, 9\}$ ;  $E(X) = 4 = \mu$

d: How spread out is the dist. of  $X$  around its mean  $\mu$ ?

$(X - \mu) \sim \text{Uniform} \{-3, -2, +5\}$   
← deviation from  $\mu$

Could try calculating  $E(X - \mu)$ , but this is 0 for any rv  $X$ , because of cancellation of  $\oplus$  and  $\ominus$  deviations; two different

easy fixes:  $E|X - \mu| \stackrel{\text{mean}}{=} \text{average absolute deviation (AAD) (MAD)}$   
↳ Gauss  $\kappa$  (Laplace)

or  $E(X - \mu)^2 \stackrel{\text{mean}}{=} \text{variance of rv } X$ .

AAD not used much; variance used constantly.

Def |  $X$  rv with finite mean  $E(X) = \mu$ ; (186)

variance of  $X = V(X) \triangleq E[(X - \mu)^2]$ .

If we  $\text{Var}(X)$  If  $E(X) = \pm\infty$  or  $E(X)$  doesn't exist,  $V(X)$  doesn't exist.

One problem with variance } The units are wrong: if  $X$  is in \$,  $V(X)$  is in \$<sup>2</sup>

Easy fix: standard deviation  $\triangleq \sqrt{V(X)} \triangleq \text{SD}(X)$  of  $X$

Consequences of these definitions

$$\textcircled{1} V(X) = E[(X - \mu)^2] = E(X^2 - 2\mu X + \mu^2)$$

$$= E(X^2) - 2\mu \underbrace{E(X)}_{\mu} + \mu^2$$

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

this is a different way to compute the variance

so  $V(X) = (\text{expectation of } X^2) - (\text{square of expectation of } X)$  (187)

Toy example

$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix}$$

$X \sim \text{Uniform}\{1, 2, 9\}$

mean  $\mu = 4$

$$E(X - \mu)^2 = \frac{1}{3}(1-4)^2 + \frac{1}{3}(2-4)^2$$

$$+ \frac{1}{3}(9-4)^2 = 12.7$$

$$SD(X) = \sqrt{12.7} = 3.6$$

This is a reasonable summary of the length of the arrows  $(= V(X))$ .

(2) For any rv  $X$ ,  $V(X) \geq 0$ ; if  $X$  is bounded,  $V(X)$  exists & is finite.

This is a consequence of Jensen's Inequality:  
 $f(x) = x^2$  is convex so  $E(X^2) \geq [E(X)]^2$ ,  
 i.e.  $V(X) = E(X^2) - [E(X)]^2 \geq 0$ .

③  $V(X) = 0 \iff P(X=c) = 1$  for some constant  $c$  (this is a trivial rv)

Notation In the same way that, by

convention,  $E(X) = \mu_X$ ,  $V(X) \equiv \sigma_X^2$

and  $SD(X) \equiv \sigma_X$  (lower-case sigma)

④  $X$  rv,  $Y = aX + b$

$\rightarrow V(Y) = a^2 V(X) = a^2 \sigma_X^2$  and

$SD(Y) = |a| \sigma_X$ . (for any constants  $a, b$ )

Special cases  $a = 1$ :  $V(X+c) = V(X)$   
 $SD(X+c) = SD(X)$

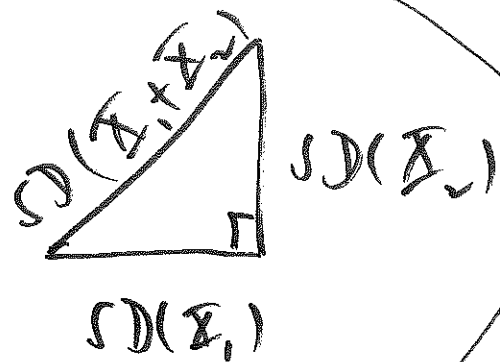
$V(aX) = a^2 V(X)$   
 $(b=0) SD(aX) = |a| SD(X)$  ⑤ If  $X_1, \dots, X_n$  are independent rv with

finite means,  $V\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n V(X_i)$ .

This is why the concept of variance (18) has endured even though the units of the variance are wrong: for

independent rvs, variance is additive,  
 whereas  $\sigma^2$  is not. name of SD: Karl Pearson (1890)  
 Special case of (5).

$X_1, X_2$  independent  $\rightarrow V(X_1 + X_2) = V(X_1) + V(X_2)$



$$\sigma = \sqrt{\sigma^2} = \sqrt{[\sigma(X_1)]^2 + [\sigma(X_2)]^2}$$

ie., SD grows like the hypotenuse of a right triangle.

Immediately,  $\max\{\sigma(X_1), \sigma(X_2)\} < \sigma(X_1 + X_2) < \sigma(X_1) + \sigma(X_2)$  (indep)

Consequence of (5)  $X_1, \dots, X_n$  independent r.v., (150)  
 $a_1, \dots, a_n, b$  constants  $\rightarrow$

$$V\left[\sum_{i=1}^n a_i X_i + b\right] = \sum_{i=1}^n a_i^2 V(X_i)$$

Example)  $X \sim \text{Binomial}(n, p)$ ; we already know that  $E(X) = np$ ; what about  $V(X)$  and  $SD(X)$ ?

Let  $S'_i = \begin{cases} 1 & \text{if success on } i^{\text{th}} \text{ trial} \\ 0 & \text{else} \end{cases}$   
for  $(i=1, \dots, n)$  and suppose as usual that

$S_1, \dots, S_n$  are IID Bernoulli( $p$ ) —

then  $X = \sum_{i=1}^n S_i$  and we can work out its variance without difficulty.

$$V(\mathbf{X}) = V\left(\sum_{i=1}^n S_i\right) \stackrel{\text{independence}}{=} \sum_{i=1}^n V(S_i) \quad \text{so } \textcircled{19/}$$

we need to work out

the variance of a Bernoulli r.v. we already know that  $E(S_j) = p$ , so if we use the formula  $V(S_j) = E(S_j^2) - [E(S_j)]^2$  we're halfway there.

Bernoulli rvs are funny:  $S_j = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } (1-p) \end{cases}$

so  $S_j^2 = \begin{cases} 1^2 = 1 & \text{with probability } p \\ 0^2 = 0 & \text{with probability } (1-p) \end{cases}$

so  $E(S_j^2) = E(S_j) = p$  and finally

$$V(S_j) = E(S_j^2) - [E(S_j)]^2 = p - p^2 = p(1-p)$$

and  $V(\bar{X}) = \sum_{i=1}^n V(S_i) = \sum_{i=1}^n p(1-p) = \boxed{n p(1-p)}$  192

and  $SD(\bar{X}) = \sqrt{n p(1-p)}$ . Example: T-S disease

$\bar{X} = (\# \text{ T-S babies in family of } n=5, \text{ both parents carriers so } p = P(\text{T-S baby}) = \frac{1}{4})$

$\sim \text{Binomial}(n, p) = \text{Binomial}(5, \frac{1}{4})$

We already worked out that  $E(\bar{X}) = np = 1.25$

Now  $SD(\bar{X}) = \sqrt{n p(1-p)} = \sqrt{5(\frac{1}{4})(\frac{3}{4})}$

It's useful to summarize  $= 0.97$   
 $= 1$

this by saying "The number of T-S babies

this couple will have will be around 1.25,

give or take about 1  $\leftarrow \sigma_{\bar{X}}$





How do you measure the spread of a distribution if the variance doesn't exist?

Example

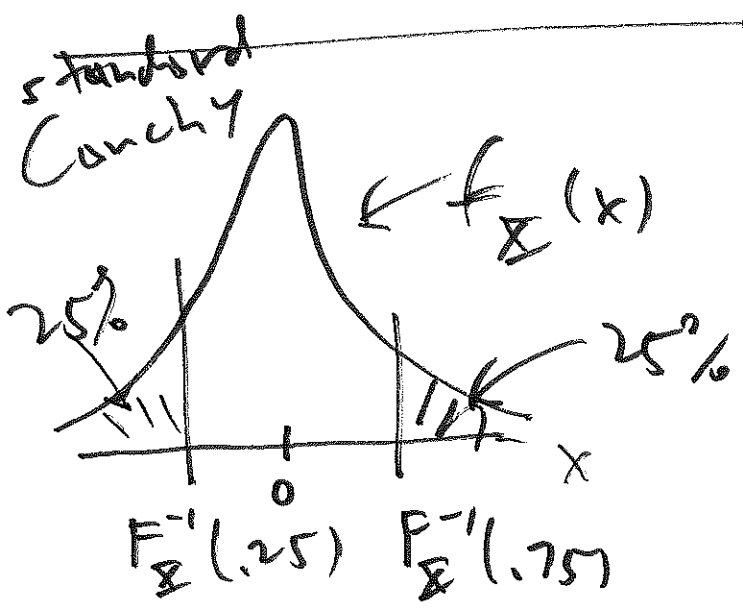
$X \sim$  (standard) Cauchy (193)

$$f_X(x) = \left\{ \begin{array}{l} \frac{1}{\pi(1+x^2)} \\ \text{for all } -\infty < x < \infty \end{array} \right\}$$

Earlier we saw that  $E(X)$  doesn't exist, so clearly  $V(X)$  doesn't exist either.

But we can use the idea of quantiles on any dist., whether its variance exists or not.

Earlier we



defined the interquartile range (IQR) as

$$\text{IQR} = \underline{\underline{F_X^{-1}(0.75) - F_X^{-1}(0.25)}}$$

standard  
Cauchy CDF is  $F_X(x) = \int_{-\infty}^x \frac{1}{\pi(1+t^2)} dt$  (194)

(arctangent)  
Here  $\tan^{-1}(x)$  is  
(calculator)  
 $= \frac{1}{2} + \frac{\tan^{-1}(x)}{\pi}$

what's called the principal

inverse of  $\tan(x)$ , varying from  $-\frac{\pi}{2}$  to

$+\frac{\pi}{2}$  as  $-\infty < x < \infty$

Need to solve

$$F_X(x) = \frac{1}{2} + \frac{\tan^{-1}(x)}{\pi} = p \text{ for } x;$$

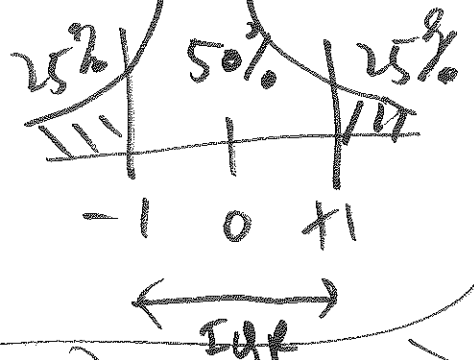
$$\text{result is } x = F_X^{-1}(p) = \tan\left(\frac{p - \frac{1}{2}}{\pi}\right),$$
$$= -\cot(p\pi)$$

So the IdR  
standard  
for the Cauchy  
distribution is

$$\text{IdR} = F_X^{-1}\left(\frac{3}{4}\right) - F_X^{-1}\left(\frac{1}{4}\right)$$
$$= \tan\left(\frac{\pi}{4}\right) - \tan\left(-\frac{\pi}{4}\right)$$
$$= 2.$$

standard Cauchy PDF

# Moments of a rv



$$E(X) = E(X^1)$$

$$V(X) = E(X^2) - [E(X^1)]^2$$

$$= E(X - \mu)^2 = E(X^2)$$

With the usual mathematical impulse to generalize:

Def.  $X$  rv,  $k$  integer  $\geq 1 \rightarrow$

$E(X^k) \triangleq$  the  $k^{\text{th}}$  moment of  $X$

of course  $E(X^k)$  may not exist, and if it does it may be infinite, but the idea is still useful.

You can show

that  $(k^{\text{th}} \text{ moment of } X \text{ exists}) \iff E(|X|^k) < \infty$

Consequences  
of the  
moment  
definition

① IF  $E(|X|^k) < \infty$  for (196)  
some integer  $k \geq 1$ , then  
 $E(|X|^j) < \infty$  for all integers  
 $j < k$ ;

in other words, if the  $k^{\text{th}}$   
moment of  $X$  exists, so do the  
 $(k-1)^{\text{st}}$ ,  $(k-2)^{\text{nd}}$ , ..., moments.

Definition

$X$  rv with expectation  $E(X) = \mu$ ,  $k$   
integer  $\geq 1 \rightarrow E[(X - \mu)^k]$  is called  
the  $k^{\text{th}}$  central moment of  $X$  or  
the  $k^{\text{th}}$  moment of  $X$  around its mean.

Clearly this idea generalizes the  
variance of  $X = E[(X - \mu)^2]$

$$\textcircled{2} \quad E[(X - \mu)^2] = E(X) - \mu = \mu - \mu = 0, \quad \textcircled{197}$$

i.e., every rv has 2<sup>nd</sup> central moment 0.

the dist. of

---

\textcircled{3} If  $X$  is symmetric around  $\mu_X$ ,

$$\text{then } E[(X - \mu)^k] = 0 \text{ for all odd}$$

integers  $k$  for which  $E[(X - \mu)^k]$  exists

---

This motivates a new definition:

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

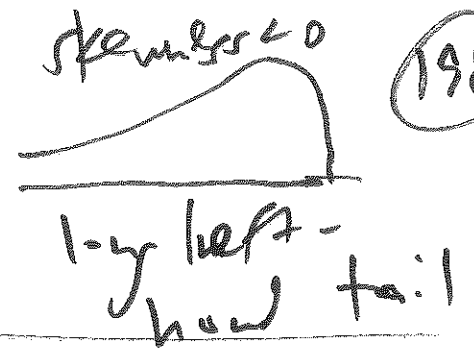
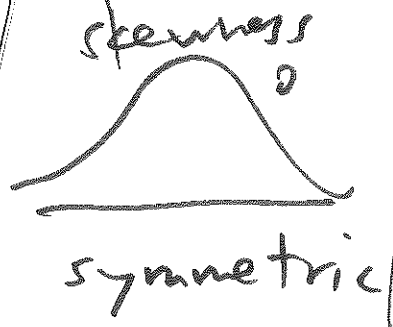
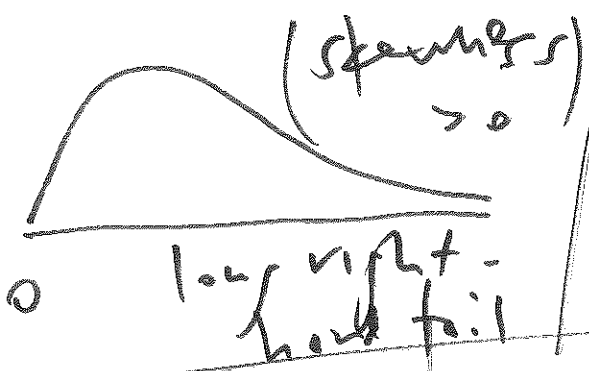
if the third moment of  $X$  exists and

is finite, then skewness ( $X$ )  $\triangleq E\left(\frac{X - \mu_X}{\sigma_X}\right)^3$ .

---

All symmetric distributions

with finite 3<sup>rd</sup> moment have skewness 0.



Moment generating functions

Def. If  $X$  rv,  $t$  a real number

$\psi_X(t) = E(e^{tX})$  is called the moment generating function of  $X$

(MGF)

The reason for this definition

Theorem If  $X$  rv with MGF  $\psi_X(t)$ , finite for all values of  $t$  in an

open interval  $(-a, b)$  around 0 ( $a > 0, b > 0$ );

then for all integers  $n > 0$ ,

$$E(X^n) = \left. \frac{d^n}{dt^n} \psi_X(t) \right|_{t=0}$$

←  $n$ th derivative of  $\psi_X$ , evaluated at  $t=0$ .

This is a handy theorem: if its premise is satisfied & the calculations are manageable, you get all the moments of  $X$  just by computing  $\psi_X(t)$  and differentiating it over & over.

Example

$X \sim \text{Exponential}(\lambda)$

( $\lambda > 0$ )

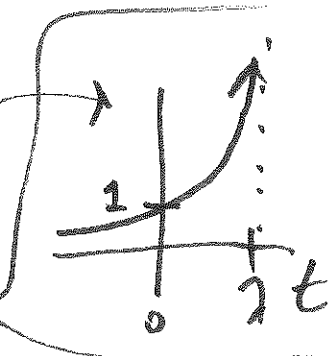
$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{else} \end{cases}$$

$$\psi_X(t) = E(e^{tX}) = \int_0^{\infty} e^{tx} \cdot \lambda e^{-\lambda x} dx = \lambda \int_0^{\infty} e^{(t-\lambda)x} dx$$

Now this integral is finite only if  $t - \lambda < 0$  is for  $t < \lambda$ , but this means (since  $\lambda > 0$ ) ~~finite~~  $- \lambda < t < \lambda$

that it's definitely finite in an open interval around 0 (eg.  $(-\lambda, \lambda)$ ).

So  $\psi(t)$  exists for  $t < \lambda$  and equals (200)

$$\psi(t) = \lambda \int_0^{\infty} e^{-(t-\lambda)x} dx = \frac{\lambda}{\lambda-t}$$


Now we just crank out the derivatives:

$$E(X) = \left. \frac{d}{dt} \frac{\lambda}{\lambda-t} \right|_{t=0} = \frac{1}{\lambda}$$

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

$$E(X^2) = \left. \frac{d^2}{dt^2} \left( \frac{\lambda}{\lambda-t} \right) \right|_{t=0} = \frac{2}{\lambda^2} = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2}$$

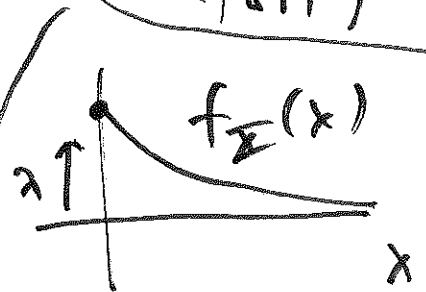
$$E(X^3) = \left. \frac{d^3}{dt^3} \left( \frac{\lambda}{\lambda-t} \right) \right|_{t=0} = \frac{6}{\lambda^3}$$

and  $SD(X) = \frac{1}{\lambda}$

$$E(X^4) = \left. \frac{d^4}{dt^4} \left( \frac{\lambda}{\lambda-t} \right) \right|_{t=0} = \frac{24}{\lambda^4}$$

positive skew (long right-hand tail)

Evidently  $E(X^k) = \frac{k!}{\lambda^k}$





Consequences  
of the  
MGF definition

①  $X$  rv with MGF  $\psi_X(t)$ , (201)

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

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

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

Example

$X$  - 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:  $\psi_{S_i}(t) = E(e^{tS_i})$

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

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

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

This was the  
Law of the  
unconscious  
Statistician

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

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

~~MGF of Binomial, continued~~

(18 Aug 17)

Since the  $X_i$  are IID,

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

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

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

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