

Simple example of this result

(X_1, X_2) joint continuous PDF (51)

$$f_{X_1, X_2}(x_1, x_2), \quad Y = a_1 X_1 + a_2 X_2 + b$$

with $a_1 \neq 0 \rightarrow Y$ continuous

with PDF $f_Y(y) = \int_{-\infty}^{\infty} f_{X_1, X_2}\left(\frac{y-b-a_2 x_2}{a_1}, x_2\right) \frac{dx_2}{|a_1|}$

Important Special case

The simplest thing you can do with two ^{or more} rvs is to add them.

This is also important in statistics, where the sample mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ plays a key role.

In the result above, take $(a_1, a_2, b) = (1, 1, 0)$ to get $Y = X_1 + X_2$.

Dist. of Y is called the convolution of the dists. of X_1 and X_2

By the above result

$$f_{\Sigma}(y) = \int_{-\infty}^{\infty} f_{\Sigma_1}(y-z) f_{\Sigma_2}(z) dz$$

A more complicated example

$$= \int_{-\infty}^{\infty} f_{\Sigma_1}(z) f_{\Sigma_2}(y-z) dz$$

is defined to be

$\Sigma_i \overset{\text{IID}}{\sim} \text{CDF } F_{\Sigma_i}, \text{ PDF } f_{\Sigma_i}$
($i=1, \dots, n$) (continuous)

$\Sigma_{(1)} \triangleq \min(\Sigma_1, \dots, \Sigma_n)$
 $\Sigma_{(n)} \triangleq \max(\Sigma_1, \dots, \Sigma_n)$ } these are examples of the order statistics of

(Take-Home test 2 problem 4) ($\Sigma_1, \dots, \Sigma_n$)

$$F_{\Sigma_{(n)}}(t) = P(\Sigma_{(n)} \leq t)$$

↓ iff (check)

$$= P(\Sigma_1 \leq t, \Sigma_2 \leq t, \dots, \Sigma_n \leq t)$$

Ⓡ IID

$$= P(\Sigma_1 \leq t) \cdot \dots \cdot P(\Sigma_n \leq t)$$

Ⓡ IID

$$= [F_{\Sigma_i}(t)]^n$$

so $Z_{(n)}$ has PDF $f_{Z_{(n)}}(t) = \frac{d}{dt} [F_{Z_0}(t)]^n$ (153)

Similarly $= n [F_{Z_0}(t)]^{n-1} f_{Z_0}(t)$

$$F_{Z_{(n)}}(t) = P(Z_{(n)} \leq t) = 1 - P(Z_{(n)} > t)$$

$$= 1 - P(X_1 > t, \dots, X_n > t) \quad \downarrow \text{if (check)}$$

$$\stackrel{\text{IID}}{=} 1 - P(X_1 > t) \dots P(X_n > t)$$

$$\stackrel{\text{IID}}{=} 1 - [1 - F_{Z_0}(t)]^n$$

so $Z_{(n)}$ has PDF $f_{Z_{(n)}}(t) = \frac{d}{dt} F_{Z_{(n)}}(t)$

$$= n [1 - F_{Z_0}(t)]^{n-1} f_{Z_0}(t)$$

Generalizing
the earlier
differentiable
& 1-1
result

Multivariate transformations 154

X_1, \dots, X_n continuous joint
dist with joint PDF $f_{\underline{X}}(\underline{x})$

support of (X_1, \dots, X_n) under $f_{\underline{X}}$

Suppose, there is a subset S of \mathbb{R}^n with

$$P[(X_1, \dots, X_n) \in S] = 1.$$

Define new vrs:

$$Y_1 = h_1(X_1, \dots, X_n)$$

\vdots

$$Y_n = h_n(X_1, \dots, X_n)$$

(note
some
arr #
of Y s)

Assume that the n
functions h_1, \dots, h_n
define a 1-1
differentiable

transformation of S onto

some subset T of \mathbb{R}^n . ← image
of h_1, \dots, h_n

Inverse
transformation:

$$x_1 = h_1^{-1}(y_1, \dots, y_n)$$

\vdots

$$x_n = h_n^{-1}(y_1, \dots, y_n)$$

Then the joint PDF $f_{\underline{Z}}(\underline{z})$ is

$$f_{\underline{Z}}(\underline{z}) = \begin{cases} f_{\underline{X}}[h_1^{-1}(\underline{z}), \dots, h_n^{-1}(\underline{z})] |J| & \text{for } (z_1, \dots, z_n) \in T \\ 0 & \text{else} \end{cases}$$

in which

J is the determinant of the matrix

$$\begin{bmatrix} \frac{\partial h_1^{-1}}{\partial z_1} & \dots & \frac{\partial h_1^{-1}}{\partial z_n} \\ \vdots & & \vdots \\ \frac{\partial h_n^{-1}}{\partial z_1} & \dots & \frac{\partial h_n^{-1}}{\partial z_n} \end{bmatrix}$$

and $| \cdot |$ is absolute value

(chain rule generalization)

J is called the Jacobian of the transformation from \underline{X} to \underline{Z} .

named after the German mathematician

Carl Gustav Jacob Jacobi (1804 - 1851)

(died of smallpox at age 46)

Looks like a generalization of the derivative of the inverse in the earlier result.

Example (X_1, X_2) joint

(continuous) PDF $f_{X_1, X_2}(x_1, x_2) = \begin{cases} 4x_1 x_2 & \text{for } 0 < x_1 < 1 \\ & 0 < x_2 < 1 \\ 0 & \text{else} \end{cases}$

(check: $\int_0^1 \int_0^1 4x_1 x_2 dx_1 dx_2$
 $= \int_0^1 4x_2 \left(\int_0^1 x_1 dx_1 \right) dx_2 = 4 \int_0^1 x_2 \left(\frac{x_1^2}{2} \Big|_0^1 \right) dx_2$
 $= 2 \int_0^1 x_2 dx_2 = 2 \left(\frac{x_2^2}{2} \Big|_0^1 \right) = 1$)

Let's work out the joint PDF of

$(Y_1, Y_2) \triangleq \left(\frac{X_1}{X_2}, X_1 \cdot X_2 \right)$

$y_1 = h_1(x_1, x_2)$
 $= \frac{x_1}{x_2}$

$y_2 = h_2(x_1, x_2) = x_1 x_2$

Inverse transform:

solve $\begin{cases} \frac{x_1}{x_2} = \gamma_1 \\ x_1 x_2 = \gamma_2 \end{cases}$ for (x_1, x_2) :

$$x_1 = h_1^{-1}(\gamma_1, \gamma_2) = \sqrt{\gamma_1 \gamma_2}$$

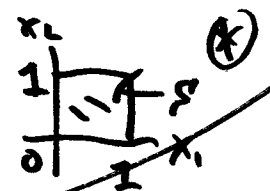
$$x_2 = h_2^{-1}(\gamma_1, \gamma_2) = \sqrt{\frac{\gamma_2}{\gamma_1}}$$

image: how does

⊗ defines 4 inequalities

$(0 < x_1 < 1, 0 < x_2 < 1)$

transform?



$$\begin{cases} x_1 > 0, x_1 < 1, \\ x_2 > 0, x_2 < 1 \end{cases}$$

So $\begin{matrix} (a) \sqrt{\gamma_1 \gamma_2} > 0, & (b) \sqrt{\gamma_1 \gamma_2} < 1 \\ (c) \sqrt{\frac{\gamma_2}{\gamma_1}} > 0, & (d) \sqrt{\frac{\gamma_2}{\gamma_1}} < 1 \end{matrix}$

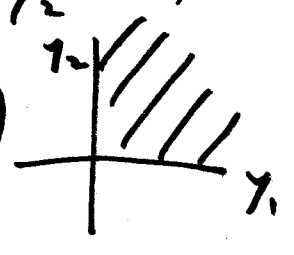
(a) equivalent to $\begin{pmatrix} \gamma_1 > 0 \\ \gamma_2 > 0 \end{pmatrix}$ or $\begin{pmatrix} \gamma_1 < 0 \\ \gamma_2 < 0 \end{pmatrix}$

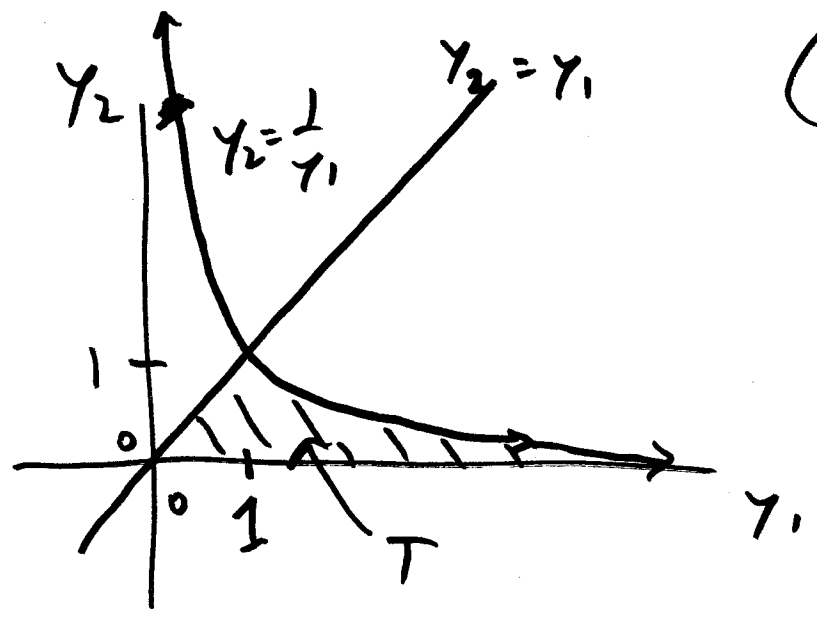
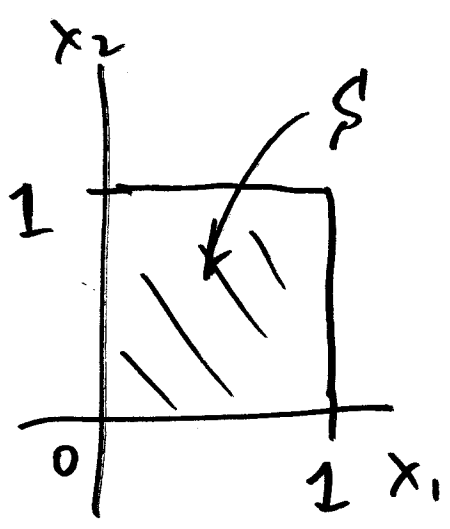
but $\gamma_1 = \frac{x_1}{x_2} > 0$ so it must be $\begin{pmatrix} \gamma_1 > 0 \\ \gamma_2 > 0 \end{pmatrix}$

(c) leads to the same thing

(b) says $\gamma_2 < \frac{1}{\gamma_1}$

(d) says $\gamma_2 < \gamma_1$





$$h_1^{-1}(y_1, y_2) = \sqrt{y_1 y_2}$$

$$h_2^{-1}(y_1, y_2) = \sqrt{\frac{y_2}{y_1}}$$

$$\text{So } \frac{d}{dy_1} h_1^{-1} = \frac{1}{2} \sqrt{\frac{y_2}{y_1}}$$

$$\frac{d}{dy_2} h_1^{-1} = \frac{1}{2} \sqrt{\frac{y_1}{y_2}}$$

$$\frac{d}{dy_1} h_2^{-1} = -\frac{1}{2} \left(\frac{y_2}{y_1^3}\right)^{\frac{1}{2}}$$

$$\frac{d}{dy_2} h_2^{-1} = \frac{1}{2} \sqrt{\frac{1}{y_1 y_2}}$$

recall
 $\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc$

$$\text{So } J = \det \begin{bmatrix} \frac{1}{2} \left(\frac{y_2}{y_1}\right)^{\frac{1}{2}} & \frac{1}{2} \left(\frac{y_1}{y_2}\right)^{\frac{1}{2}} \\ -\frac{1}{2} \left(\frac{y_2}{y_1^3}\right)^{\frac{1}{2}} & \frac{1}{2} \left(\frac{1}{y_1 y_2}\right)^{\frac{1}{2}} \end{bmatrix} = \frac{1}{2y_1}$$

and (since $y_1 > 0$) $|J| = \frac{1}{2y_1}$

To finish the calculation, in the

$$\text{PDF of } \underline{X}, f_{\underline{X}}(\underline{x}) = \begin{cases} 4x_1 x_2 & (0 < x_1 < 1) \\ & (0 < x_2 < 1) \\ 0 & \text{else} \end{cases}$$

substitute $x_1 = \sqrt{y_1 y_2}$, $x_2 = \sqrt{\frac{y_2}{y_1}}$
 and bring in the Jacobian:

$$f_{\underline{Y}}(\underline{y}) = f_{\underline{X}}[h_1^{-1}(\underline{y}), h_2^{-1}(\underline{y})] |J|$$

$$= 4 \left(\sqrt{y_1 y_2} \right) \left(\sqrt{\frac{y_2}{y_1}} \right) \frac{1}{2y_1}$$

$$= \begin{cases} 2 \frac{y_2}{y_1} & \text{for } (y_1, y_2) \in T \\ 0 & \text{else} \end{cases}$$

A useful
trick

start with (X_1, X_2) joint 160
dist.; suppose you're interested
only in the dist. of $Y_1 = h_1(X_1, X_2)$.
Then one way to compute this dist. is
with the following ³ steps.

Step 1: Find

another $Y_2 = h_2(X_1, X_2)$ such that
the transformation $(X_1, X_2) \rightarrow (Y_1, Y_2)$ is
1-1 with a differentiable inverse transformation
& the calculations are straightforward.

Step 2 work out the joint dist. of

(Y_1, Y_2) . Step 3 Integrate Y_2 out of

the joint dist. (i.e., marginalize over
 Y_2) to get the marginal dist. of Y_1 .

Example of
4 \mathcal{I}_2 that
would not work

$$\mathcal{I}_1 = 2\mathcal{X}_1$$

$$\mathcal{I}_2 = 3\mathcal{X}_1 = \frac{3}{2}\mathcal{I}_1$$

(161)

Here \mathcal{I}_2 is linearly dependent on \mathcal{I}_1 , so the rank of the ^(2x2) Jacobian matrix is only 1 and its determinant is therefore 0.

~~Earlier~~

Example,
continued

$(\mathcal{X}_1, \mathcal{X}_2)$ have

joint (continuous) PDF

$$f_{\mathcal{X}_1, \mathcal{X}_2}(x_1, x_2) = \begin{cases} 4x_1x_2 & 0 < x_1 < 1 \\ & 0 < x_2 < 1 \\ & 0 & \text{else} \end{cases}$$

Earlier

we found

$$\text{that with } (\mathcal{Y}_1, \mathcal{Y}_2) = \begin{pmatrix} \mathcal{X}_1 \\ \mathcal{X}_2 \\ \mathcal{X}_1, \mathcal{X}_2 \end{pmatrix}$$

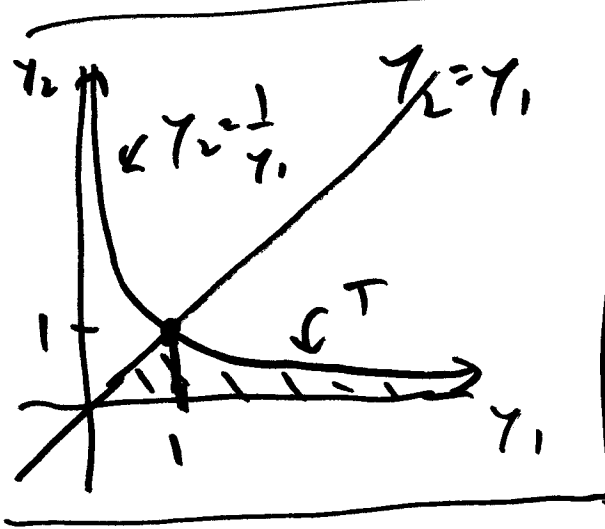
the transformed

PDF was

$$f_{\mathcal{Y}_1, \mathcal{Y}_2}(y_1, y_2) = \begin{cases} \frac{2y_2}{y_1} & \text{for } (y_1, y_2) \in T \\ 0 & \text{else} \end{cases}$$

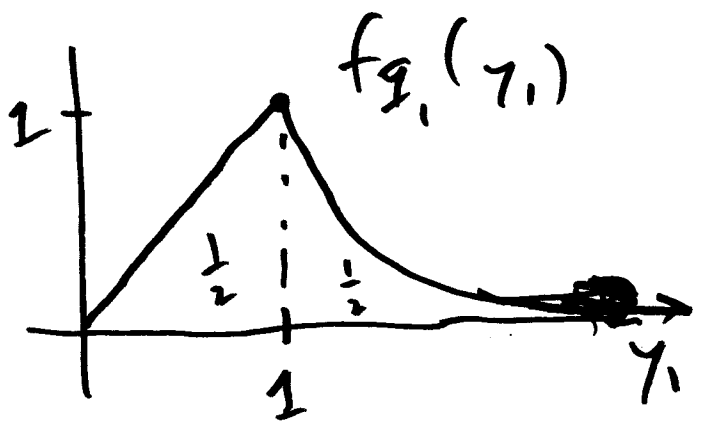
where $T = \{(y_1, y_2) : y_1 > 0, y_2 < \min(y_1, \frac{1}{y_1})\}$.

Suppose you were only really interested
 (marginal)
 in the dist. of $Z_1 = \frac{X_1}{X_2}$; then all you have
 to do is integrate Z_2 out of the joint dist.



For $y_1 > 0$, the allowable
 region for y_2 is in two
 parts: for $0 < y_1 < 1, 0 < y_2 < y_1$
 and for $y_1 > 1, 0 < y_2 < \frac{1}{y_1}$

$$\text{So } f_{Z_1}(y_1) = \begin{cases} \int_0^{y_1} 2\left(\frac{y_2}{y_1}\right) dy_2 = y_1 & \text{for } 0 < y_1 < 1 \\ \int_0^{1/y_1} 2\left(\frac{y_2}{y_1}\right) dy_2 = y_1^{-3} & \text{for } y_1 > 1 \end{cases}$$



weird PDF: not
~~not~~ differentiable
 at $y_1 = 1$

Useful consequence of Jacobian story

$\underline{X} = (X_1, \dots, X_n)$ continuous with joint PDF $f_{X_1, \dots, X_n}(x_1, \dots, x_n)$

$\underline{Y} = (Y_1, \dots, Y_n)$ is a linear transformation of \underline{X}

$\underline{Y}^T = A \cdot \underline{X}^T$ where A is an invertible (full-rank) matrix

matrix.

Then the PDF of \underline{Y} is

$$f_{\underline{Y}}(\underline{y}) = \frac{f_{\underline{X}}(A^{-1} \underline{y})}{|\det A|}$$

Example

$$Y_1 = X_1 + X_2$$

$$Y_2 = X_1 - X_2$$

$$A = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$\det A = -2 = ad - bc$$

$$|\det A| = 2$$

(recall that)

$$A^{-1} = \frac{1}{-2} \begin{pmatrix} -1 & -1 \\ -1 & 1 \end{pmatrix} = \frac{1}{2} A$$

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \rightarrow A^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

Expectation,
Variance,
Covariance,
Correlation

we showed

Disch. 4

Example: T-5 (184)
disorder (continued)

Earlier we worked out the
discrete dist. of the rv

$I = (\# \text{ of T-5 babies in family})$
(of 5, both parents carriers)

that $(I) \sim \text{Binomial}(n, p)$ with $\begin{cases} n=5 \\ p=\frac{1}{4} \end{cases}$

y	$P(I=y)$
0	$\binom{5}{0} \left(\frac{1}{4}\right)^0 \left(\frac{3}{4}\right)^5 = 0.2373$
1	$\binom{5}{1} \left(\frac{1}{4}\right)^1 \left(\frac{3}{4}\right)^4 = 0.3955$
2	$\binom{5}{2} \left(\frac{1}{4}\right)^2 \left(\frac{3}{4}\right)^3 = 0.2637$
3	$\binom{5}{3} \left(\frac{1}{4}\right)^3 \left(\frac{3}{4}\right)^2 = 0.0879$
4	$\binom{5}{4} \left(\frac{1}{4}\right)^4 \left(\frac{3}{4}\right)^1 = 0.0146$
5	$\binom{5}{5} \left(\frac{1}{4}\right)^5 \left(\frac{3}{4}\right)^0 = 0.0010$
	1.0000

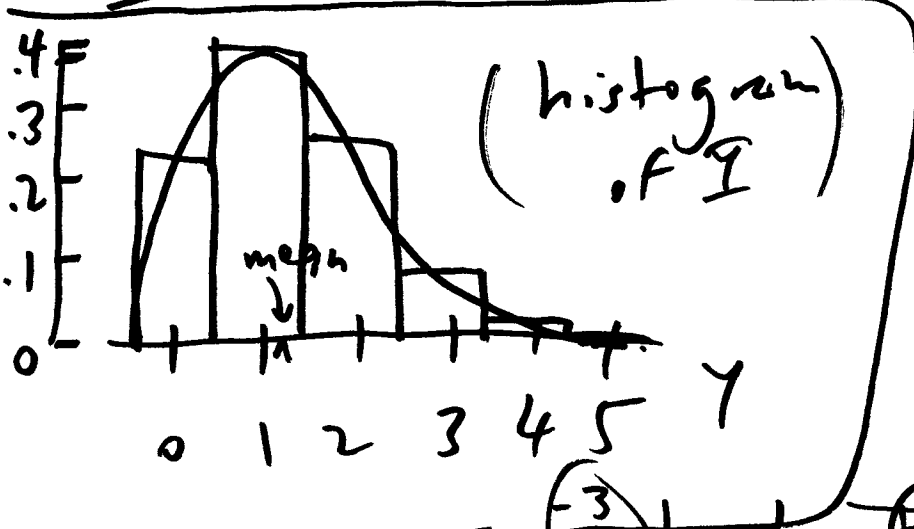
$$P(I=y) = \begin{cases} \binom{n}{y} p^y (1-p)^{n-y} & y=0,1,\dots,n \\ 0 & \text{else} \end{cases}$$

Q: About how
many T-5 babies
should these parents
expect to have?

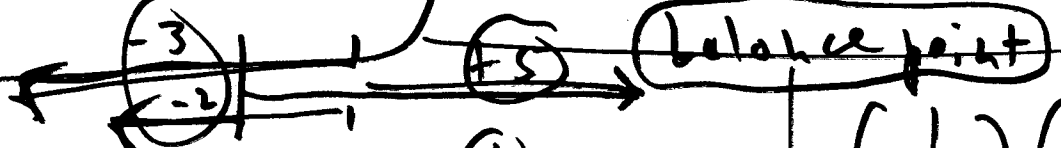
(center of dist.?
of I)

A₁ Most likely outcome is 1 T-S body (165)
 (mode of the dist. of \mathcal{Y})

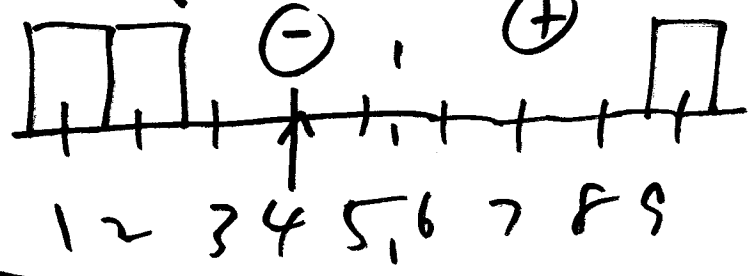
A₂ (physics idea)



let's work out the center of mass of the distribution



toy example



$$\begin{pmatrix} 1 \\ 2 \\ \textcircled{9} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

outlier

let's find the place c where the histogram balances: where (the sum of forces exerted by the histogram bars to the left of c) equals (the sum of forces to the right):

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \rightarrow \begin{pmatrix} y_1 - c \\ \vdots \\ y_n - c \end{pmatrix}$$

want sum = 0

$$\sum_{i=1}^n (y_i - c) = 0 = \left(\sum_{i=1}^n y_i \right) - nc = 0$$

A₃ Median of the dist. of \mathcal{I} (here that's also 1)

$$\sum_{i=1}^n y_i - nc = 0 \iff$$

$$c = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} = \text{the sample mean of the (sample) dataset}$$

here $\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix}$ mean $\bar{y} = 4$

Here each value of \mathcal{I} occurred only once:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\bar{y} = \sum_{i=1}^n \left(\frac{1}{n}\right) y_i \quad \text{Def.}$$

If some values are more probable than others, the generalization of $\left(\frac{1}{n}\right)$ weight on each y value would be to weight each y by its probability $P(\mathcal{I} = y)$.

A rv is bounded if all of its possible values are finite.

Def.

let \mathcal{I} be a bounded discrete rv with PF $\frac{1}{n}$

$$f_{\mathcal{I}}(y) = P(\mathcal{I} = y). \text{ The}$$

mean or expected value or expectation of \mathcal{I} ,

is $E(Z) \triangleq \sum_{\text{all } y} y P(Z=y) = \sum_{\text{all } y} y f_Z(y)$ (16)

T-s example

$$E(Z) = (0)(.2373) + (1)(.3955)$$

$$+ \dots + (5)(.0012) = 1.2500000$$

suspiciously round #

Symbolically if $Z \sim \text{Binomial}(n, p)$

then $E(Z) = \sum_{y=0}^n y \binom{n}{y} p^y (1-p)^{n-y}$

Wol

$$= \sum_{y=1}^n y \binom{n}{y} p^y (1-p)^{n-y}$$

since summand is 0 for $y=0$

Wolfram alpha

$$= \sum_{y=1}^n y \frac{n!}{y!(n-y)!} p^y (1-p)^{n-y}$$

cancel y against $y \cdot (y-1)!$

$$= \sum_{y=1}^n \frac{n \cdot (n-1)!}{y(y-1)! \cdot (n-1-(y-1))!} p \cdot p^{y-1} (1-p)^{n-y}$$

$$= np \sum_{y=1}^n \frac{(n-1)!}{(y-1)!(n-y)!} p^{y-1} (1-p)^{n-1-(y-1)}$$

$\binom{n-1}{y-1}$

This assumes that $n > 1$; proof for $n = 1$ is on the next page

$$= np \sum_{y=1}^n \binom{n-1}{y-1} p^{y-1} (1-p)^{n-1-(y-1)} \quad (168)$$

$$= np \sum_{i=0}^{n-1} \binom{n-1}{i} p^i (1-p)^{n-1-i} \quad \left(\begin{array}{l} \text{substitute} \\ i=y-1 \end{array} \right)$$

So: if

Binomial($n-1, p$)
dist.

$Z \sim \text{Binomial}(n, p)$

for $n \geq 1$, $E(Z) = np$

this = 1
because binomial
probabilities add
up to 1

When $n=1$, Binomial($1, p$) = Bernoulli(p).

In this case $E(Z) = 0 \cdot P(Z=0) + 1 \cdot P(Z=1)$

So: for all
 $n \geq 1$ (integer)

$$= 0 \cdot (1-p) + 1 \cdot p = p$$

$$= np \text{ with } n=1$$

and $0 < p < 1$, $Z \sim \text{Binomial}(n, p) \rightarrow E(Z) = np$.

T-S example) $(n=5, p=\frac{1}{4})$ $E(X) = \frac{5}{4} = 1.25$ (169) ✓

If discrete X is unbounded, the expectation of X may not exist, ^{either} because

$$\sum_{x < 0} x f_X(x) = -\infty \quad (\text{and/or} \quad \sum_{x \geq 0} x f_X(x) = +\infty)$$

or the distribution "puts too much mass

near $\pm\infty$ "

Def. X discrete rv with

$\sum_{x < 0} f_X(x)$; consider $\sum_{x < 0} x f_X(x)$ and

$\sum_{x \geq 0} x f_X(x)$. If both sums are infinite,

$E(X)$ is undefined (or does not exist);

if at least one sum is finite, then

$$E(X) = \sum_{\text{all } x} x f_X(x) \text{ exists} \quad \left(\begin{array}{l} \text{it} \\ \text{may} \\ \text{still} \\ \text{be} \\ \text{infinite} \end{array} \right)$$

To create a discrete rv whose mean doesn't exist, you have to play a careful game, because $\sum_{\text{all } x} f_{\mathbb{I}}(x)$ has to be finite (it has to equal 1) but $\sum_{\text{some } x} x f_{\mathbb{I}}(x)$ has

to be infinite.

Example

The harmonic

series $(\frac{1}{1} + \frac{1}{2} + \frac{1}{3} + \dots) = \sum_{x=1}^{\infty} \frac{1}{x}$ was known

to the ancient Greeks, because ^(integers) the wavelengths of the overtones of a vibrating string are $\frac{1}{2}, \frac{1}{3}, \dots$ of the fundamental wavelength

of the string. The fact that $\sum_{x=1}^{\infty} \frac{1}{x} = \infty$

(i.e., the harmonic series diverges) was first ^{French} shown in the 1300s (!) by the philosopher Nicole Oresme (~1320-1382).

It's clear from this divergence that (171)
you can't create a rv X with $P^m F$

$P(X=x) = \frac{c}{x}$, $x=1, 2, \dots$, because the
probability ^{would} sum to $+\infty$.

$$\text{But } P(X=x) = \frac{c}{x^2}$$

or $P(X=x) = \frac{c}{x(x+1)}$ turn out to work;

for example, $\sum_{x=1}^{\infty} \frac{1}{x^2} = \frac{\pi^2}{6}$ (!) and, even

more conveniently, $\sum_{x=1}^{\infty} \frac{1}{x(x+1)} = 1$.

If we use this to construct two pathological
discrete distributions, to show what can
go wrong with the idea of expectation.

Example 1

$$f_X(x) = \begin{cases} \frac{1}{x(x+1)} & x=1, 2, \dots \\ 0 & \text{else} \end{cases}$$

$$E(X) = \sum_{x=1}^{\infty} x \cdot \frac{1}{x(x+1)} = \sum_{x=1}^{\infty} \frac{1}{x+1} = +\infty \quad (172)$$

so $E(X)$ exists, it's just infinite.

Example 2

$$f_X(x) = \begin{cases} \frac{1}{2|x|(1+|x|)} & x = \pm 1, \pm 2, \dots \\ 0 & \text{else} \end{cases}$$

we already know that $\sum_{\text{all } x} f_X(x) = 1$, so X is a well-defined rv; but $\sum_{x=-\infty}^{\infty} x \cdot \frac{1}{2|x|(1+|x|)} =$

and $\sum_{x=1}^{\infty} x \cdot \frac{1}{2x(x+1)} = +\infty$, so $E(X)$

does not exist.

we won't work with pathological rv, mostly.

Expectation
for continuous
rvs

Def. X bounded
continuous rv

with PDF $f_X(x) \rightarrow E(X) \stackrel{\Delta}{=} \int_{-\infty}^{\infty} x f_X(x) dx$ (173)

Example) $X \sim \text{Exponential}(\lambda)$ ($\lambda > 0$):

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

$$\text{So } E(X) = \int_0^{\infty} \lambda x e^{-\lambda x} dx \stackrel{\text{integrate by parts}}{=} \frac{1}{\lambda}$$

For this reason, many people parameterize the exponential distribution differently:

Alternative definition

$X \sim \text{Exponential}(\eta)$ ($\eta > 0$)

$$\rightarrow f_X(x) = \begin{cases} \frac{1}{\eta} e^{-\frac{x}{\eta}} & x > 0 \\ 0 & \text{else} \end{cases}$$

with this parameterization

you can see that $E(X) = \eta$ (easier to interpret).

Nevertheless, to avoid confusion with (174)
DS, I'll stick with $\lambda e^{-\lambda x}$.

If continuous
rv \mathcal{I} is unbounded, a bit of care is once
again required to define $E(\mathcal{I})$.

Def.

\mathcal{I} continuous rv with PDF $f_{\mathcal{I}}(y)$; consider

$$\int_{-\infty}^0 y f_{\mathcal{I}}(y) dy \quad \text{and} \quad \int_0^{\infty} y f_{\mathcal{I}}(y) dy.$$

If both integrals are infinite, $E(\mathcal{I})$ is
undefined (or does not exist); if

at least one of these integrals is

finite, $E(\mathcal{I}) = \int_{\mathbb{R}} y f_{\mathcal{I}}(y) dy$ exists

(but it may still be infinite).

Example 2 A dist. that does arise in 175
practical statistical applications is
the Cauchy distribution (attributed
to Augustin-Louis Cauchy (1789-1857)
a French mathematician who wrote 800
25 textbooks & 800 research articles in a 52-year period (15/year,
but actually first studied carefully by
Poisson in 1824).

$$f_{\mathcal{C}}(y) = \frac{1}{\pi(1+y^2)} \quad (-\infty < y < \infty)$$

is the (standard) Cauchy distribution.

It does integrate to 1, but $\int_0^{\infty} \frac{y}{\pi(1+y^2)} dy = \infty$

and $\int_{-\infty}^0 \frac{y}{\pi(1+y^2)} dy = -\infty$, so $E(\mathcal{C})$ does not exist,

because its tails go to 0 extremely slowly.

This is because for large γ , $\frac{\gamma}{1+\gamma^2} \approx \frac{1}{\gamma}$

and $\int_c^\infty \frac{1}{\gamma} d\gamma = +\infty$, the continuous analogue of the harmonic series

(ch7 c70)

analogue of the harmonic series

Expectation of a function of a rv

~~rv~~ \mathbb{R} continuous rv with PDF $f_{\mathbb{X}}(x)$, $\mathbb{I} \triangleq h(\mathbb{X})$.

Method 1

work out PDF $f_{\mathbb{I}}(\gamma)$;

$$E(\mathbb{I}) = \int_{\mathbb{R}} \gamma f_{\mathbb{I}}(\gamma) d\gamma$$

(if this exists)

Method 2 (faster)

$$E(\mathbb{I}) = \int_{\mathbb{R}} h(x) f_{\mathbb{X}}(x) dx$$

Discrete version:

$$E[h(\mathbb{X})] = \sum_{\text{all } x} h(x) f_{\mathbb{X}}(x)$$

↑
discrete

DS (and some other people) call Method 2 (177) ^(L0745)
the Law of the Unconscious Statistician,

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

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

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

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

~~(177)~~

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 $E(\underline{X}) \leq b$.

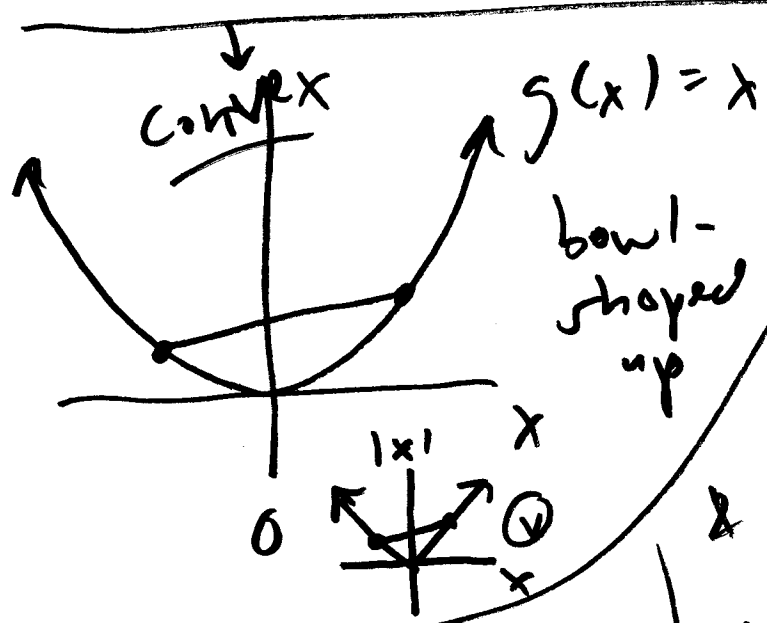
③ If $\underline{X}_1, \dots, \underline{X}_n$ are n rvs, each with finite $E(\underline{X}_i)$, then $E\left(\sum_{i=1}^n \underline{X}_i\right) = \sum_{i=1}^n E(\underline{X}_i)$,

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 .

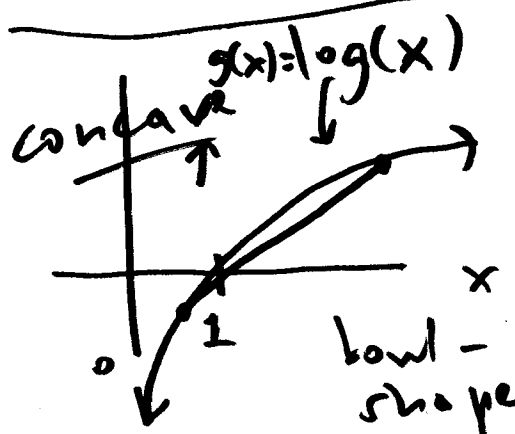
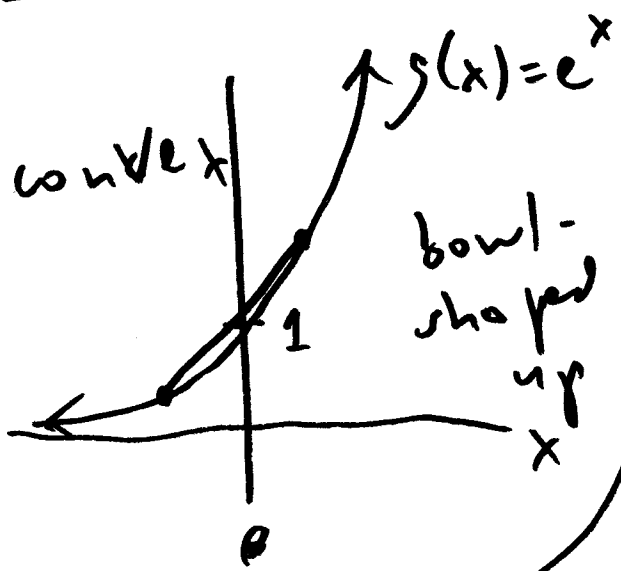
Def. A function $g: \mathbb{R}^n \rightarrow \mathbb{R}$ (this means that $g(\underline{x}) = z$ is convex
 $\begin{array}{c} \nwarrow \quad \nearrow \\ \text{real \#s} \\ \nwarrow \quad \nearrow \\ (x_1, \dots, x_n) \end{array}$)

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

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



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 entirely above the function except at the endpoints.



g is concave if

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

concave

Def. The expectation of a random vector

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

(a) g convex, \tilde{X} random vector with finite

$E(\tilde{X}) \rightarrow E[g(\tilde{X})] \geq g[E(\tilde{X})]$

Jensen's Inequality

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

(attributed to Johan Jensen (1859-1925),

Denish mathematician & engineer)
(Йенсен)

Applications of (3)

Suppose that $X_1, \dots, X_n \stackrel{IID}{\sim}$ 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_j are
independent

X_1, \dots, X_n independent rv, each with
finite $E(X_j) \rightarrow$ (181)

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

Contrast this with 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

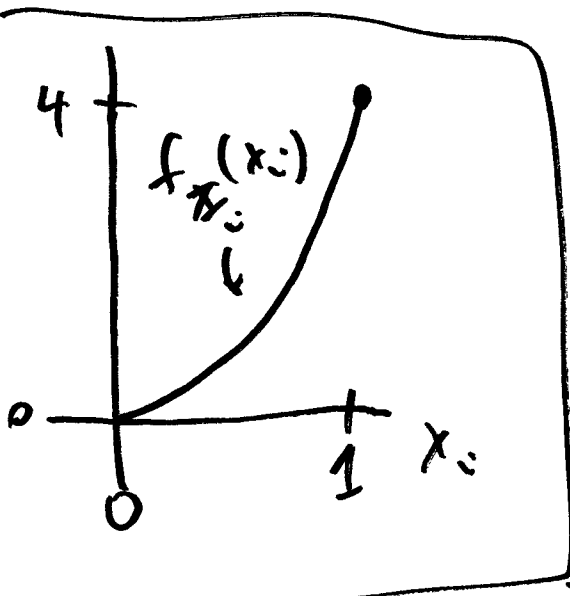
its job if you filtered the water twice
instead of once?

X_1 = proportion of bad stuff removed in the 1st filtering (184)

X_2 = proportion removed in 2nd filtering of what was left from 1st 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 Z = proportion of original bad stuff remaining after 2 filtrations = $(1-X_1)(1-X_2)$

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

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 $E(Z) = \mu^2$.

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

so 20% of bad stuff expected to be removed in 1st filtering; $E(I) = \mu^2 = 0.04$, so expect only 4% of bad stuff to remain after 2 filterings.

(b) Suppose

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

(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 (a)

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 184

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

(1st subgoal)

process to take; $E(X) = ?$

geometric distribution

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 geometric series

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

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

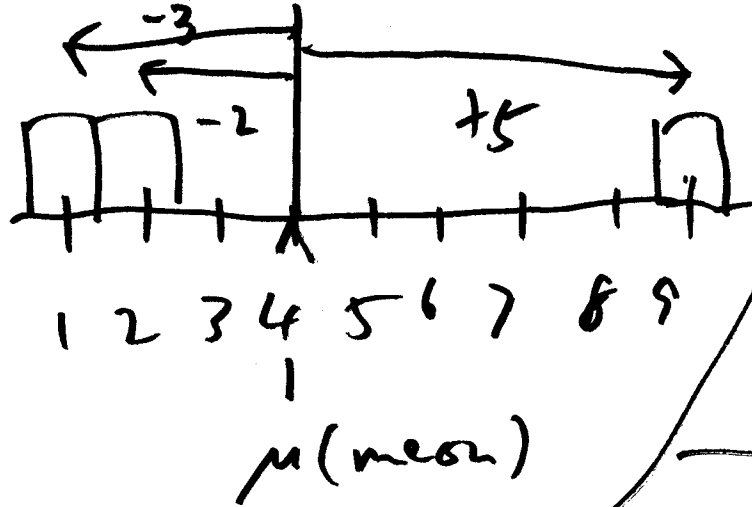
(inverse relation \checkmark)

If I'm terrible ($p = .01$)

I expect to succeed on

the $\frac{1}{.01} = 100$ th throw.

Variance
and
standard
deviation



(185)

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

$$\text{mean } \bar{x} = \mu$$

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

Q: How spread out is the dist. of X around its mean μ ?

← deviation from μ

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

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: $E|X - \mu| \stackrel{\Delta}{=} \text{average absolute deviation (AAD) (MAD)}$

← Gauss (1809) ← Laplace (1785)

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

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
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 messy: if
 X is in \$, $V(X)$ is in \$².

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

of X \leftarrow (K. Pearson [name]) (~1890)

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

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

Toy example $\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix}$ $X \sim \text{Uniform } \{1, 2, 9\}$
 mean $\mu = 4$ $E(X - \mu)^2 = \frac{1}{3}(1-4)^2 + \frac{1}{3}(2-4)^2 + \frac{1}{3}(9-4)^2 = 12.7$
 $\text{SD}(X) = \sqrt{12.7} = 3.6$ $(= V(X))$

This is a reasonable summary of the length of the arrows

(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. $V(X) = E(X^2) - [E(X)]^2 \geq 0$. concave

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, $Y = aX + b$

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

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

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

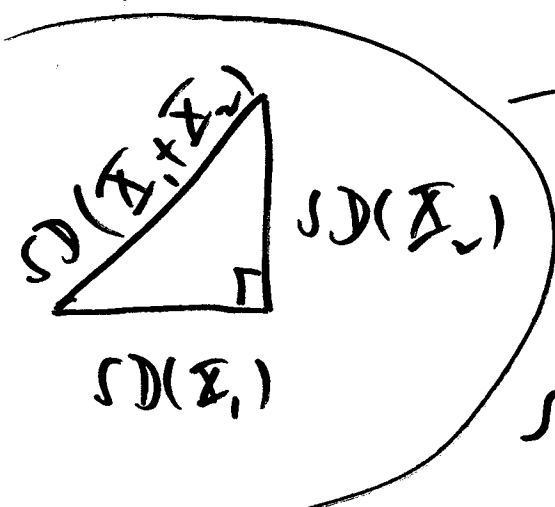
$V(aX) = a^2 V(X)$
 $(b=0) SD(aX) = |a| SD(X)$ ⑤ If X_1, \dots, X_n
 are independent rv with

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 (189) has endured even though the units of the variance are wrong: for

independent rvs, variance is additive
 when we hear \rightarrow SD is not correct units. (16 Aug 19) name of SD: Karl Pearson (1890)
 Special case of (5).

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



SD

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

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

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

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

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

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

Let $S_i = \begin{cases} 1 & \text{if success on } i^{\text{th}} \text{ success} \\ 0 & \text{failure trial} \\ 0 & \text{else} \end{cases}$
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.

$$V(\mathbf{X}) = V\left(\sum_{i=1}^n S_i\right) \stackrel{\text{independence}}{=} \sum_{i=1}^n V(S_i) \quad \text{So } \textcircled{191}$$

we need to work out

the variance of a Bernoulli rv. 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 $V(\underline{X}) = \sum_{i=1}^n V(S_i) = \sum_{i=1}^n p(1-p) = \boxed{np(1-p)}$ (192)

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

\underline{X} = (# T-S babies in family of $n=5$, 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 $E(\underline{X}) = np = 1.25$

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

It's useful to summarize $\underline{X} = 0.97$
 \downarrow
 $= 1$

this by saying "The number of T-S babies

this couple will have will be around 1.25,

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

\uparrow
 $\mu_{\underline{X}}$

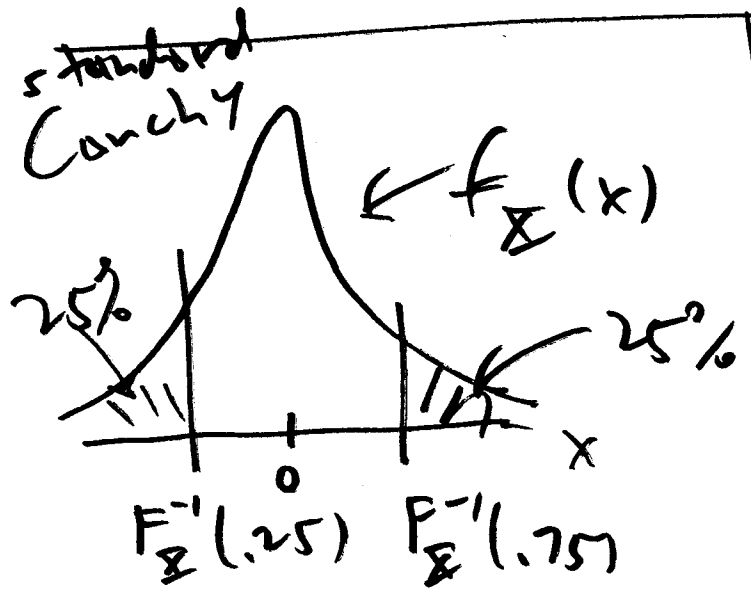
How do you measure the spread of a distribution if the variance doesn't exist?

Example $X \sim$ (standard) Cauchy

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

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

$$\text{IQR} = \underline{F_X^{-1}(.75)} - F_X^{-1}(.25)$$

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 $-\infty < x < \infty$

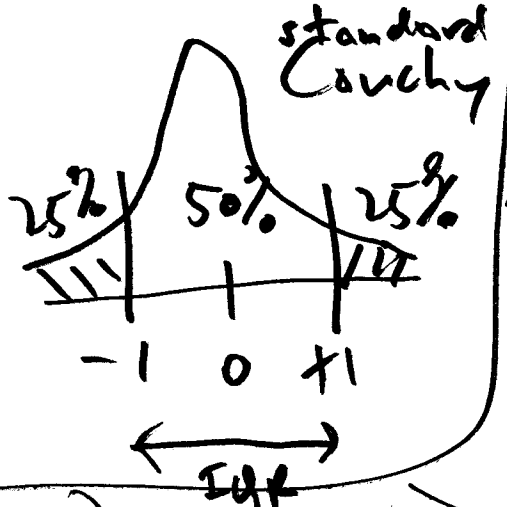
Need to solve

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

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

So the IdR
standard
for the Cauchy
distribution is

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



Moments of a rv

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

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

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

with the

usual mathematical impulse to generalize:

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

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

Definition

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

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

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

ie, every rv has 2nd central moment 0.

~~the dist. of~~

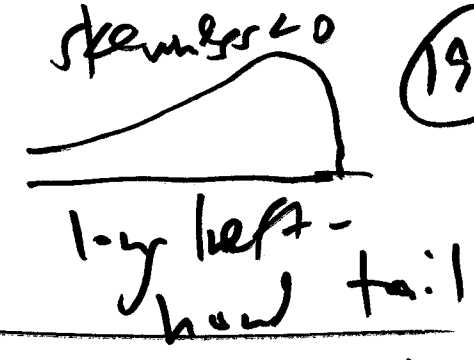
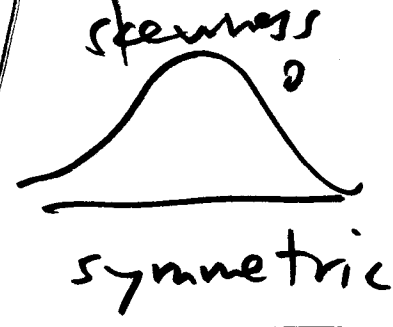
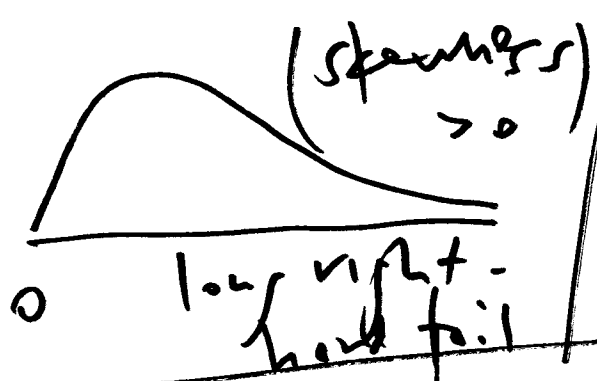
$\textcircled{3}$ If X is symmetric around μ_X ,
then $E[(X - \mu)^k] = 0$ for all odd
integers k for which $E[(X - \mu)^k]$ exists

This motivates a new definition:

Def X rv with mean μ_X , SD σ_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 3rd moment have skewness 0.

converting X
to standard units



Moment generating functions

Def. | X rv, t a real number

$\psi(t) \triangleq E(e^{tX})$ is called the moment generating function of X
 (MGF)

The reason for this definition

Theorem | 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 ψ_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.

(16 May 19)

Example

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

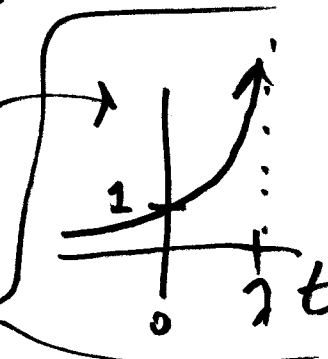
$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0 & \text{else} \end{cases}$$

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

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

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

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

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


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

