

Multivariate distributions

So far we've looked at (127) one or two rvs at a time; easy to generalize to a finite number of rv  $Z_1, \dots, Z_n$ ,  $n$  positive finite integer.

Def. The joint CDF of  $n$  rvs  $Z_1, \dots, Z_n$  is the function  $F_{Z_1, \dots, Z_n}(y_1, \dots, y_n)$  specified by  $F_{Z_1, \dots, Z_n}(y_1, \dots, y_n) = P(Z_1 \leq y_1, \dots, Z_n \leq y_n)$

More compact to use vector notation:  $\underline{Z} = (Z_1, \dots, Z_n)$ ,  $\underline{y} = (y_1, \dots, y_n)$

$F_{\underline{Z}}(\underline{y}) = P(Z_1 \leq y_1, \dots, Z_n \leq y_n)$   $\underline{Z}$  is

said to be a random vector taking values in  $\mathbb{R}^n$ .

Def.  $n$  rv  $(Z_1, \dots, Z_n) \stackrel{\sim}{=} Z$  have a discrete joint distribution if the random vector  $Z$  can only take on a finite or countably infinite # of possible values  $(z_1, \dots, z_n) \in \mathbb{R}^n$ .

The joint PF (probability <sup>mass</sup> function) of  $Z$

is  $f_{Z_1, \dots, Z_n}(z_1, \dots, z_n) = P(Z_1 = z_1, \dots, Z_n = z_n)$

or equivalently  $f_Z(z) = P(Z = z)$ .

Example  $n$  patients in treatment group of a randomized clinical trial;  $B_i = \begin{cases} 1 & \text{if patient } i \text{ has a good outcome} \\ 0 & \text{else} \end{cases}$

If nothing else is known about the patients (e.g., age, disease burden at start of trial, ...) it would be reasonable to model the  $B_i$  as IID Bernoulli  $(\theta)$  <sup>same</sup> success probability.

$\underline{B} = (B_1, \dots, B_n)$ ;  $\underline{b} = (b_1, \dots, b_n)$ ;  $\underline{B}$  has a 129  
 discrete joint distribution  $f_{\underline{B}}(\underline{b}) = P(B_1 = b_1, \dots, B_n = b_n)$ .  
PF ↖ 100 ↗

If  $\theta$  were known you could use  $f_{\underline{B}}(\underline{b})$  to predict the dataset before it arrives: by the IID assumption  $P(B_1 = b_1, \dots, B_n = b_n | \theta) = P(B_1 = b_1) \dots P(B_n = b_n)$

Recall that  $P(B_i = b_i | \theta) = \theta^{b_i} (1-\theta)^{1-b_i}$  for  $b_i = 0, 1$  so

$$\begin{aligned}
 f_{\underline{B}}(\underline{b}) &= \prod_{i=1}^n \theta^{b_i} (1-\theta)^{1-b_i} = \theta^{\sum_{i=1}^n b_i} (1-\theta)^{n - \sum_{i=1}^n b_i} \\
 &= \theta^s (1-\theta)^{n-s}
 \end{aligned}$$

Def.  $n$  rv  $Z_1, \dots, Z_n$  have a continuous joint distribution if you can find a function  $f_{\underline{Z}}$  on  $\mathbb{R}^n$  such that for every (non-weird) subset  $\bullet \subset \mathbb{R}^n$  with  $s = \sum_{i=1}^n b_i$ .

$$P[(Z_1, \dots, Z_n) \in G] = \int \dots \int_G f_{Z_1, \dots, Z_n}(z_1, \dots, z_n) dz_1 \dots dz_n$$

$f_{\underline{Z}}(z)$  is the joint PDF (probability density function) of  $\underline{Z}$ .

more compactly

$$P(\underline{Z} \in G) = \int \dots \int_G f_{\underline{Z}}(z) dz$$

Consequences of this def.

① If the joint dist. of  $\underline{Z}$  is continuous,

then  $f_{\underline{Z}}(z) = \frac{\partial^n}{\partial z_1 \dots \partial z_n} F_{\underline{Z}}(z)$

Mixed discrete/continuous

with n rv

random vectors behave just as they do with 2 rv.

more realistically,  $\theta$  would

Example clinical trial (continued)

be unknown, and you can think about the

joint dist. of  $(\underline{B}, \theta) = (B_1, \dots, B_n, \theta)$ , (131)  
 in which the  $B_i$  are discrete and  $0 < \theta < 1$  is  
 continuous.

### Marginal distributions

If you know the joint PDF  $f_{\underline{Z}}$  of  $\underline{Z}$ , you  
 can work out the marginal distribution of  
any subset of  $(Z_1, \dots, Z_n)$  by integrating  
 ~~$f_{\underline{Z}}$~~   $f_{\underline{Z}}(\underline{z})$  over the elements of  $(Z_1, \dots, Z_n)$   
 that are not in the subset.

### Example

$$\underline{Z} = (Z_1, Z_2, Z_3, Z_4)$$

$$f_{Z_1}(z_1) = \iiint f_{\underline{Z}}(\underline{z}) dz_2 dz_3 dz_4$$

$$f_{Z_2, Z_3}(z_2, z_3) = \iint f_{\underline{Z}}(\underline{z}) dz_1 dz_4 \quad \text{and so on.}$$

Similarly, you can work out a marginal  
 CDF by sending the other components

to  $\infty$ : for example

(132)

$$F_{\underline{Z}}(\underline{z}) = P(\underline{Z} \leq \underline{z}) = P(Z_1 \leq z_1, Z_2 < \infty, \dots, Z_n < \infty)$$

$$= \lim_{z_2 \rightarrow \infty, \dots, z_n \rightarrow \infty} F_{\underline{Z}}(\underline{z})$$

Definition

$n$  rvs  $Z_1, \dots, Z_n$  are independent if  
non-weight

for any  $n$  sets  $A_1, \dots, A_n$  of real numbers

$$P(Z_1 \in A_1, \dots, Z_n \in A_n) = \prod_{i=1}^n P(Z_i \in A_i)$$

Immediate  
consequences

①  $Z_1, \dots, Z_n$  independent iff

$$F_{\underline{Z}}(\underline{z}) = \prod_{i=1}^n F_{Z_i}(z_i)$$

②  $Z_1, \dots, Z_n$

independent iff  $f_{\underline{Z}}(\underline{z}) = \prod_{i=1}^n f_{Z_i}(z_i)$

(133)

Def. Starting with a univariate  $P_n^M$  or (133)

PDF  $f_{\Sigma_i}(y_i)$ ,  $n$  rvs  $(\Sigma_1, \dots, \Sigma_n)$  form a random sample of size  $n$  from  $f_{\Sigma_i}$  if the  $\Sigma_i$  are

independent and all of them have marginal  $P_n^M$  or PDF  $f_{\Sigma_i}$   $\leftrightarrow$  i.e., if the  $\Sigma_i$  are an independent identically distributed (IID)

sample from  $f_{\Sigma_i}$

Example

deer at usc:  
some have a disease  
(chronic wasting disease)

population  
all deer living within usc boundary

9 May 2019

Sample  
the observed deer

disease?  
↑  
 $N = ?$   
( $\approx 800$ )  
↓  
 $\begin{bmatrix} 1s \\ & \\ 0s \end{bmatrix}$

mean  $\theta = ?$   
(unknown)

~~IID~~

disease?  
 $\begin{bmatrix} 1 \\ 2 \\ \vdots \\ n \end{bmatrix} \begin{bmatrix} 1s \\ & \\ 0s \end{bmatrix}$

mean  $\bar{y} = \hat{\theta}$   
"y-bar"  
↑

$1 = y$   
 $0 = N$   
 $n = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

estimate of  
"theta-hat"

Short hand for the diagram:

$$(Y_i | \theta) \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$$

(end example for now)

Definition (134)

Start with random vector

$\underline{X} = (X_1, \dots, X_n)$ ; partition it into 2

subvectors  $\underline{X} = (\underline{Y}, \underline{Z})$ ,  $\underline{Y} = (Y_1, \dots, Y_k)$   
 $1 \leq k \leq n-1$

$$\underline{Z} = (Z_1, \dots, Z_{n-k})$$

Then for every point

$\underline{z}$  for which  $f_{\underline{Z}}(\underline{z}) > 0$ , the conditional

distribution of  $\underline{Y}$  given  $\underline{Z}$  is

$$f_{\underline{Y} | \underline{Z}}(\underline{y} | \underline{z}) = \frac{f_{\underline{Y}, \underline{Z}}(\underline{y}, \underline{z})}{f_{\underline{Z}}(\underline{z})}, \quad \underline{y} \in \mathbb{R}^k$$

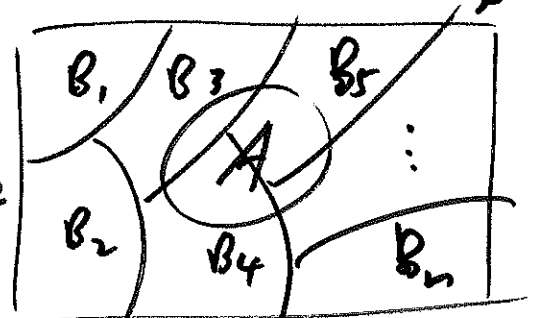
from which

$$f_{\underline{Y}, \underline{Z}}(\underline{y}, \underline{z}) = f_{\underline{Z}}(\underline{z}) f_{\underline{Y} | \underline{Z}}(\underline{y} | \underline{z}).$$



Multivariate  
law of total  
probability

You'll recall that if 135  
A is an event & you're



trying to compute  $P(A)$  & it's hard, one idea is to find another aspect of the world  $B$  upon which  $A$  depends, such that the events  $B_1, \dots, B_n$  form a partition;

$$\text{then } P(A) = \sum_{i=1}^n P(A \cap B_i) = \sum_{i=1}^n P(B_i) P(A|B_i)$$

This has an analogue with continuous r.v.s.

using the notation in the definition of conditional distributions

$$f_{\underline{X}}(\underline{x}) = \int \dots \int_{\mathbb{R}^{n-k}} \underbrace{f_{\underline{X}}(\underline{z})}_{\text{like } B_i} \underbrace{f_{\underline{X}|\underline{Z}}(\underline{x}|\underline{z})}_{\text{like } P(A|B_i)} d\underline{z}$$

Multivariate  
Bayes's  
Theorem

using the same notation, (136)

$$f_{\tilde{z}|\tilde{y}}(\tilde{z}|\tilde{y}) = \frac{f_{\tilde{z}}(\tilde{z}) f_{\tilde{y}|\tilde{z}}(\tilde{y}|\tilde{z})}{f_{\tilde{y}}(\tilde{y})}$$

(posterior info)      (prior info)      (likelihood info)

unknown data

The usual application of this in statistics is as follows.

(normalizing constant)  $f_{\tilde{y}}(\tilde{y})$

Def.  $\tilde{z}$  a random vector with multivariate distribution  $f_{\tilde{z}}(\tilde{z})$ ; then random variables  $X_1, \dots, X_n$  are conditionally independent given  $\tilde{z}$  if for all  $\tilde{z}$  with  $f_{\tilde{z}}(\tilde{z}) > 0$ ,

$$f_{\tilde{X}|\tilde{z}}(\mathbf{x}|\tilde{z}) = \prod_{i=1}^n f_{X_i|\tilde{z}}(x_i|\tilde{z}).$$

Earlier  
example,  
revisited

Remember the machine with (137)  
a  $\theta$  dial that can make IID  
coin tosses with  $P(\text{heads}) = \theta$ !

earlier

We agreed that, if  $\theta$  is unknown to you,

① the results of the coin tosses  $\mathcal{I}_1, \mathcal{I}_2, \dots$

are dependent, because there is useful

information in any subset of them for

predicting any other subset, but ② the

$\mathcal{I}_i$  become conditionally independent

given  $\theta$ , because once you know  $\theta$

there's no longer any useful information

in the  $\mathcal{I}_i$  to predict other  $\mathcal{I}_i$ .

this is why - in both the clinical trial, example & the (nuts & bolts) example - we

model the data values  $Y_i$  as  $(Y_i | \theta) \stackrel{\text{conditionally}}{\sim} \text{Bernoulli}(\theta)$ .

Functions of a rv

Case 1: discrete

$X$  discrete rv with  $P_f = f_X(x)$ ;  $Y = h(X)$  for some function

$h$  defined on {possible values of  $X$ }. Then

$$f_Y(y) = P(Y=y) = P(h(X)=y)$$

$$= \sum_{\{X: h(X)=y\}} f_X(x)$$

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

The median of this distribution is 5;  
 $Y = |X - 5| = h(X)$  keeps track of how far  $X$  is from the median.

$Y$	$X$ such that $X+Y=Y$	$P(Z=Y)$
0	5	1/9
1	4 or 6	2/9
2	3 or 7	2/9
3	2 or 8	2/9
4	1 or 9	2/9
		<hr/> 1

Case 2: Continuous

---

$Z$  continuous  
or with PDF

$f_Z(x)$ ;  
 $Z = h(X)$   
as before.

The CDF  $F_Z(y)$  can be worked out as follows:

$$F_Z(y) = P(Z \leq y) = P[h(X) \leq y]$$

$$= \int_{\{x: h(x) \leq y\}} f_X(x) dx$$

and if  $Z$  is also continuous

$$f_Z(y) = \frac{d}{dy} F_Z(y)$$

(at every point  $y$  where  $F_Z$  is differentiable).

Example)  $\lambda$  = rate at which customers served in a queue at the bank

Natural to model  $\lambda$  as continuous, (also,  $\lambda > 0$ ) with CDF  $F_\lambda$ .

Turns out that the average waiting time is  $\bar{W} = \frac{1}{\lambda} = h(\lambda)$ . You can get the PDF of  $\bar{W}$

- in 2 steps:
- ① work out CDF of  $\bar{W}$
  - ② differentiate with respect to  $y$
- ① (for  $y > 0$ )

$$\begin{aligned}
 F_{\bar{W}}(y) &= P(\bar{W} \leq y) = P\left[h(\lambda) \leq y\right] \\
 &= P\left(\frac{1}{\lambda} \leq y\right) = P\left(\lambda \geq \frac{1}{y}\right) \quad \text{since } \lambda \text{ is continuous} \\
 &= 1 - P\left(\lambda < \frac{1}{y}\right) = 1 - P\left(\lambda \leq \frac{1}{y}\right) \\
 &= 1 - F_\lambda\left(\frac{1}{y}\right) \quad \text{and now}
 \end{aligned}$$

$$f_Z(y) = \frac{d}{dy} F_Z(y) = \frac{d}{dy} \left( 1 - F_X\left(\frac{1}{y}\right) \right) \quad (14)$$

chain rule

$$= -f_X\left(\frac{1}{y}\right) \left(-y^{-2}\right) = \frac{f_X\left(\frac{1}{y}\right)}{y^2}$$

Example

$X \sim \text{Uniform}[-1, +1]$  (14 Aug 17)  
(continuous)

$$Z = X^2$$

find PDF of  $Z$

First

note that  $Z$ 's possible values are  $[0, 1]$ .

for  $0 < y < 1$

$$\textcircled{1} F_Z(y) = P(Z \leq y) = P(X^2 \leq y)$$

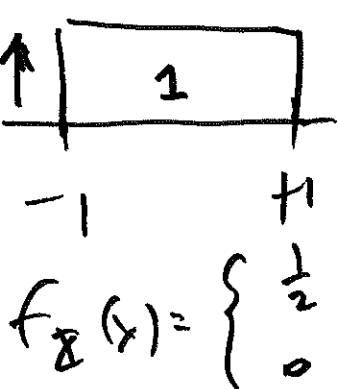
$$= P(-\sqrt{y} \leq X \leq \sqrt{y}) = \int_{-\sqrt{y}}^{\sqrt{y}} f_X(x) dx$$

$$= \frac{1}{2} x \Big|_{-\sqrt{y}}^{\sqrt{y}}$$

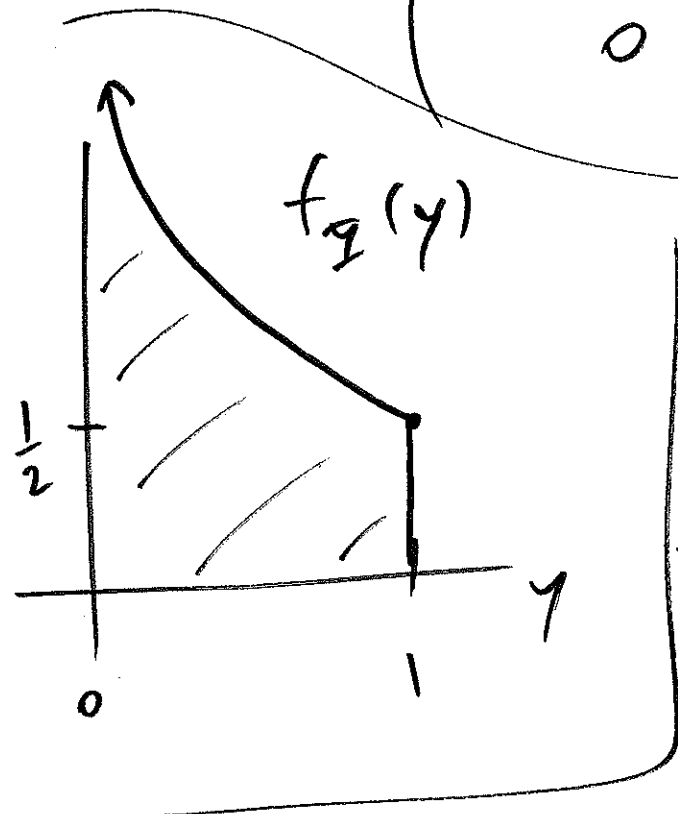
$$= \sqrt{y}$$

② Thus

$$f_Z(y) = \frac{d}{dy} F_Z(y)$$



$$\text{So } f_Y(y) = \begin{cases} \frac{1}{2\sqrt{y}} & \text{for } 0 < y < 1 \\ 0 & \text{else} \end{cases}$$



This density is unbounded at 0 (!). Every theorem

$X$  continuous rv with pdf  $f_X(x)$ ,

$$Y = aX + b \quad (a \neq 0) \quad \text{(linear transformation)}$$

$$\rightarrow f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right)$$

Interesting and useful fact

$X$  continuous with CDF  $F_X(x)$ ; what's the distribution of  $Y = F_X(X)$ ?



$$F_Z(\gamma) = P(Z \leq \gamma) = P\left[F_X(X) \leq \gamma\right] \quad (143)$$

$$\left. \begin{array}{l} \text{for} \\ 0 < \gamma < 1 \end{array} \right\} = P\left[X \leq F_X^{-1}(\gamma)\right] = F_X\left[F_X^{-1}(\gamma)\right] = \gamma$$

But the dist. with  $F_Z(\gamma) = \gamma$  for  $0 < \gamma < 1$  is the Uniform  $(0, 1)$  distribution (!)

Probability Integral Transform	with $F_X$
	$X$ continuous, CDF, $Z = F_X(X)$
	$\rightarrow Z \sim \text{Uniform}(0, 1)$ or $[ ]$

why is this useful?

Converse is also true:  
 $Z \sim \text{Uniform}[0, 1]$ ,  $F_X^{-1}$

continuous CDF with quantile function

$$F_X^{-1} \rightarrow X = F_X^{-1}(Z) \sim F_X$$

This is the practical basis for the generation of many forms of pseudo-random numbers. (144)

It turns out to be easy to generate pseudo-uniform  $(0, 1)$  values; therefore if you want to generate pseudo-random #s from a distribution with CDF  $F_X$  and  $F_X^{-1}$  is easy & fast to compute,

Algorithm

$U_1, \dots, U_n \stackrel{\text{IID}}{\sim} \text{Uniform}(0, 1)$

(Quiz 6)  $\rightarrow F_X^{-1}(U_1), \dots, F_X^{-1}(U_n) \stackrel{\text{IID}}{\sim} F_X$

Earlier Example revisited

If  $X \sim \text{Exponential}(\lambda)$ , its

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

Earlier we saw that  $F_X(x) = \begin{cases} 0 & x \leq 0 \\ 1 - e^{-\lambda x} & x > 0 \end{cases}$  (p. 91)

and  $F^{-1}$   

$$F^{-1}(p) = \frac{-\log(1-p)}{\lambda}$$
 $(0 < p < 1)$   
 (R demo)

Now  
 you can see  
 immediately

that if  $U \sim \text{Uniform}(0, 1)$  so is  $(1-U)$ ,  
 so to generate IID Exponential<sup>(2)</sup> rv you  
 just compute  $-\frac{1}{\lambda} \log U_i$ ,  $U_i \sim \text{Uniform}(0, 1)$   
 (rexp) R ~~rexp~~

why do  
 people  
 want/need  
 pseudo-  
 random  
 numbers?

Some stochastic (probabilistic)  
 models of real-world phenomena  
 are too complicated to fully  
 characterize mathematically  
 in closed form; one highly

useful method in such situations is  
 (computer-based)  
 to conduct a simulation study driven  
 by pseudo-random numbers.

Bedrock method  
 in data science  
 today.

The method used above for working out (146)  
the distribution of  $\bar{Y} = \frac{1}{X}$  can be  
generalized, or follows.

Some functions  $h(\bar{Y})$

are nice, in that they are both differentiable

and one-to-one (invertible)

Calculus  
reminder

← real-valued

If  $h(x)$  is differentiable and one-to-one (1-1)

for  $x$  in the open interval  $(a, b)$ , then

$h$  is either monotonically increasing or

decreasing, and  $h$  is also continuous,

so it transforms the interval  $(a, b)$  to

another open interval  $h[(a, b)] = (\alpha, \beta)$

called the image of  $(a, b)$  under  $h$ .

Since  $h$  is invertible, it makes sense

to talk about  $y = h(x) \Leftrightarrow x = h^{-1}(y)$ . (147)

**Theorem:**  $X$  continuous rv with PDF  $f_X(x)$  and for which  $P(a < X < b) = 1$ ;  $\Sigma = h(X)$ , with  $h$  differentiable and 1-1 for  $a < x < b$ ;  $(\alpha, \beta)$  image of  $(a, b)$  under  $h$ ;  $h^{-1}(y)$  inverse function of  $h(x)$  for  $\alpha < y < \beta$

of  $\Sigma$  is  $f_{\Sigma}(y) = \begin{cases} f_X[h^{-1}(y)] \left| \frac{dh^{-1}(y)}{dy} \right| & \text{for } \alpha < y < \beta \\ 0 & \text{else} \end{cases}$

→ PDF (chain rule)

Every short-hand way to remember this: "Multiply" both sides

$$y = h(x)$$

$$x = h^{-1}(y)$$

$$by |dy| \text{ to get } f_{\Sigma}(y) |dy| = f_X(x) |dx|$$

$\bar{X} = h(\bar{X}) = \frac{1}{\bar{X}}$ : average waiting time in the bank queue

Earlier example, revisited

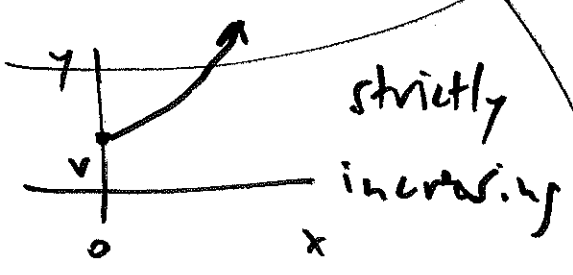
Here  $y = h(x) = \frac{1}{x}$  so  $x = h^{-1}(y) = \frac{1}{y}$

and  $\frac{d}{dy} \frac{1}{y} = -\frac{1}{y^2}$ ; thus  $f_{\Sigma}(y) = \frac{f_{\Sigma}(\frac{1}{y})}{y^2}$  as before

Example / At time 0,  $v$  organisms introduced into large tank of water with nutrients;  $\Sigma$  = rate of growth. Under one model that's realistic in some circumstances, at time  $t$  the predicted population size would be  $\Sigma = v e^{\Sigma t}$  (exponential growth).

$\Sigma$  unknown, modeled with  $f_{\Sigma}(x) = \begin{cases} 3(1-x)^2 & 0 < x < 1 \\ 0 & \text{else!} \end{cases}$

$y = h(x) = v e^{xt}$



$x=0 \rightarrow y=v$

$x=1 \rightarrow y = v e^{t}$  image



$$\frac{y}{v} = e^{xt} \rightarrow \log\left(\frac{y}{v}\right) = xt \rightarrow x = h^{-1}(y) = \frac{1}{t} \log\left(\frac{y}{v}\right) \quad (49)$$

$$\frac{d}{dy} \frac{1}{t} \log\left(\frac{y}{v}\right) = \frac{1}{t} \left(\frac{y}{v}\right)^{-1} \cdot \frac{1}{v} = \frac{1}{ty} \quad \text{Thus}$$

$$f_{\mathbb{I}}(y) = \begin{cases} \frac{3 \left[1 - \frac{1}{t} \log\left(\frac{y}{v}\right)\right]^2}{ty} & v < y < ve^t \\ 0 & \text{else} \end{cases}$$

(9 May 19)

Functions  
of 2 or  
more rvs

Case 1:  
discrete

with joint  $\prod_{i=1}^m f_{\mathbb{I}_i}(x_i)$ ;  
discrete joint dist.  
n rvs  $\mathbb{I}_1, \dots, \mathbb{I}_n$

$$\text{define } \left\{ \begin{array}{l} \mathbb{I}_1 = h_1(\mathbb{I}_1, \dots, \mathbb{I}_n) \\ \vdots \\ \mathbb{I}_m = h_m(\mathbb{I}_1, \dots, \mathbb{I}_n) \end{array} \right\} \quad (m \geq 1)$$

↑  
real-valued

$$(h_j : \mathbb{R}^n \rightarrow \mathbb{R})$$