

Follow the logic detailed on Dr. ④<sup>46</sup>  
pp. 49-50  
to obtain

$$P\left(\bigcup_{i=1}^n A_i\right) = \frac{1}{1!} - \frac{1}{2!} + \frac{1}{3!} - \dots + (-1)^{n+1} \frac{1}{n!}$$

calculator result:

$$\lim_{n \rightarrow \infty} \left[ \sum_{i=1}^n \frac{(-1)^{i+1}}{i!} \right] = 1 - \frac{1}{e} \approx 0.63$$

This sum "approaches its limit quickly"; already with  $n=7$  you have the first 4 significant figures: 0.6321

(29 July 16)

Conditional probability

Note that Kolmogorov's

probability axioms

defined the function

$P_k(A)$ , where  $A$  is a set in the

collection  $\mathcal{C}$  of subsets of the sample space  $\mathbb{S}$  in which nothing weird can occur; in other words,  $P_k(A)$  is a function of a single argument  $A$ .

To include the extremely useful idea of conditional probability in his

setup, Kolmogorov has to define it

using  $P_k$ . Definition Given any two

events  $A, B$  in  $\mathcal{C}$ , the conditional probability of  $A$  given  $B$  is

$$P(A|B) = \begin{cases} \frac{P(A \cap B)}{P(B)} & \text{if } P(B) > 0 \\ \text{undefined} & \text{if } P(B) = 0 \end{cases}$$

There are other foundational theories (48) of probability - one by the Italian mathematician and actuary Bruno de Finetti (def) (1906-1985) and another by the American physicists Richard T. Cox (1898-1991) and Edwin T. Jaynes (1922 - 1998) (CJ) - in which the probability function  $P_{\text{def}}(A|B)$  or  $P_{\text{CJ}}(A|B)$  has 2 inputs, not 1, so that conditional probability is the primitive concept, not unconditional probability or with Kolmogorov's  $P_k(A)$ . def and CJ

were responding to the reality that (49)  
in practice, all probabilities are conditional  
on background Assumptions, Information  
and Judgments (AIJ)

Example (Tay-Sachs)

we actually computed not

$P(\text{at least 1 TS baby})$  but

$P(\text{at least 1 TS baby} \mid \text{family of 5, } \text{mother and } \text{father both carriers})$

This impulse, to be explicit about your  
AIJ, is Bayesian; Kolmogorov worked  
in the frequentist paradigm; in this  
course, focusing on  $P_k(B)$ , we need to  
remember that it should really be  $P_k(B_{\text{AIJ}})$ .

(50)

Consequences  
of the  
conditional  
probability  
definition  
(theorems)

①  $A, B$  events in  $\mathcal{C}$ :

if  $P(B) > 0$  then

$$P(A \cap B) = P(B) P(A|B)$$

and if  $P(A) > 0$

then  $P(A \cap B) = P(A) P(B|A)$ .

② Direct generalization: if  $A_1, \dots, A_n$  are events with  $P(A_1 \cap \dots \cap A_{n-1}) > 0$ ;

then

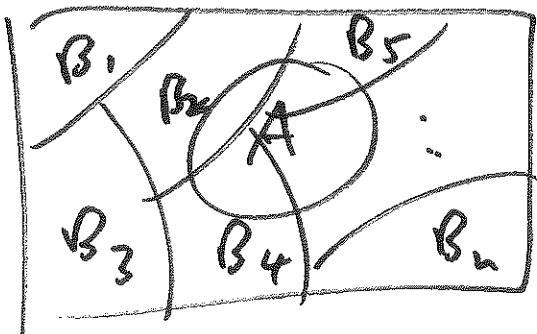
chain rule for  $\cap = \text{and}$

$$P(A_1 \cap \dots \cap A_n) = P(A_1) P(A_2 | A_1) P(A_3 | A_1 \cap A_2) \dots P(A_n | A_1 \cap \dots \cap A_{n-1})$$

Definition

$S$  sample space; if you can find events  $B_1, \dots, B_k$  in  $\mathcal{C}$

such that the  $B_j$  are disjoint and (5) exhaustive ( $\bigcup_{i=1}^n B_i = \Omega$ ), then you have found a partition  $(B_1, \dots, B_k)$ .



③ If  $(B_1, \dots, B_k)$  is a partition of  $\Omega$

with  $P(B_j) > 0$  for all  $j = 1, \dots, k$ , then for any event  $A$  in  $C$

$$P(A) = \sum_{j=1}^k P(B_j) P(A|B_j) -$$

This is the law of Total Probability

LTP

(51.1)

When is the LT p useful? } You're trying to compute  $P(A)$  and you find it hard to compute directly. If you can find some aspect  $B$  of the world satisfying 2 properties -

- ①  $B$  defines a partition  $\{B_1, \dots, B_k\}$  of  $S$  and ②  $A$  depends on  $B$  in such a way that the conditional probabilities  $P(A | B_j)$  are easier to compute than  $P(A)$  itself - then you can work out  $P(A)$  indirectly:  $P(A) = \sum_{j=1}^k \underbrace{P(B_j)}_{P(A \cap B_j)} P(A | B_j)$ .

Extension of the LTP (4) Assuming all conditional probabilities are defined (52)

is what follows, if  $C$  is in  $\mathcal{C}$  then

$$P(A|C) = \sum_{j=1}^k P(B_j|C) P(A|B_j \cap C).$$

Definition Events  $A, B$  are independent if

$$P(A \cap B) = P(A) \cdot P(B),$$

which ( $\Leftrightarrow$  by  $\Leftrightarrow P(A) > 0, P(B) > 0$ )

is equivalent to

$$P(A|B) = P(A)$$

$$\text{and } P(B|A) = P(B).$$

Consequences  
of the  
definition  
of independence

① If  $A$  and  $B$  are (53)  
independent, then so are  
 $A$  and  $B^c$ ,  $A^c$  and  $B$ ,  
and  $A^c$  and  $B^c$ .

② Extension of the definition to  
more than 2 events:

Definition:

Given events  $A_1, \dots, A_k$ , they are  
(mutually) independent if, for  
every subset  $A_{i_1}, \dots, A_{i_j}$  of  $(A_1, \dots, A_k)$   
( $j = 2, \dots, k$ ),

$$P(A_{i_1} \cap \dots \cap A_{i_j}) = P(A_{i_1}) \cdot \dots \cdot P(A_{i_j})$$

(54)

Interpretation of independence

$A, B$  independent  $\iff$  information about  $A$

doesn't change the chances associated with  $B$ , and vice versa.

Definition

Another <sup>useful</sup> extension of independence

Events  $\{A_1, \dots, A_k\}$  are conditionally independent given event  $B$  if for every subset  $\{A_{i_1}, \dots, A_{i_j}\}$  of  $\{A_1, \dots, A_k\}$  ( $j = 2, \dots, k$ )

$$P(A_{i_1} \cap \dots \cap A_{i_j} | B) = \prod_{l=1}^j P(A_{i_l} | B)$$

$\downarrow$  product

Example] Suppose that there is a machine that can take an ordinary coin and produce IID tosses of the coin with  $P(H) = \theta$ , and  $\theta$  can be set to any value in  $[0, 1]$  with a dial on the machine's control panel.

Someone sets the dial to a  $\theta$  value that's unknown to you and starts producing coin tosses  $I_1, I_2, \dots$ . Suppose the first 10 tosses come out HTHTHTTHTH (7 H, 3 T).

Q: Is there information in these first 10 tosses that helps you to predict  $I_{11}$ ?

A: Yes, definitely: it looks like  $\theta$  is around  $\frac{7}{10}$ , so you would predict  $I_{11} = H$ . Thus  $I_{11}$  depends on  $I_1 \dots I_{10}$  probabilistically. Now, suppose instead that you watched the guy with the machine set the dial to  $\theta = 0.81$ , so that  $\theta$  is now known to you. The next 10 tosses come out  $HHHHTHTHHHHH$  ( $8H, 2T$ ). Q: Is there information in these 10 tosses that helps you to predict the next toss?

A: No; you know that  $\theta = 0.81$ , so there's no information in any of the  $I_v$ .

that helps you to predict any of (57)  
 the other  $\bar{I}_j$ .  
 Given  $\theta$ , the  $\bar{I}_j$  are indep.

Thus the  $\bar{I}_j$  are  
unconditionally dependent but  
conditionally independent given  $\theta$ .

Bayes's Theorem for events Suppose that the events  $B_1, \dots, B_k$  partition the sample space in such a way that  $P(B_j) > 0$  for all  $j = 1, \dots, k$ . If  $A$  is an event with  $P(A) > 0$ , then for each  $i = 1, \dots, k$

$$P(B_i | A) = \frac{P(B_i) P(A | B_i)}{P(A)}$$

and, by the LTP, this is (58)

$$P(B_i | A) = \frac{P(B_i) P(A | B_i)}{\sum_{j=1}^k P(B_j) P(A | B_j)}$$

How this theorem is used in Bayesian

statistics

The  $B_i$  represent unknown

states of the world: They're all  
possible —  $P(B_i) > 0$  — and only one

of them is true, but you don't know  
which one.  $(A)$  represents data:

information that will help you identify  
the most probable  $B_i$ .

Before the dataset A arrives, you have background information about the plausibility of the  $B_i$  that you can represent with prior probabilities  $P(B_i)$ .

After the dataset A arrives, you can use Bayes' Theorem to update your prior probabilities to posterior probabilities  $P(B_i | A)$ .

The probabilities  $P(A | B_i)$  represent how likely the dataset A would be if  $B_i$  were the actual unknown state; this is often called likelihood information.

$P(A)$  does not depend on the  $B_i$ ,  
and can therefore be regarded as a  
normalizing constant, put into

Bayes's theorem to make all the

$P(B_i | A)$  add up to 1. But

$$P(B_i | A) = \frac{P(B_i) P(A | B_i)}{P(A)}$$

is interpreted as

$$\begin{aligned} (\text{posterior information}) &= \frac{(\text{prior information}) \cdot (\text{likelihood information})}{(\text{normalizing constant})}. \end{aligned}$$

# Random variables and their distributions (6)

Example : Tay-Sach's Disease

§ 2

$T = T\text{-S baby}$   
 $N = \text{not}$

NNNNNN	0	# of T-S babies = $\bar{Y}$
TNNNNN	1	Given a sample
NTNNNN		Definition
NNTNNN		Space $S'$ for an experiment $E$ ,
NNNTNN		a (real-valued) random variable
NNNNNT		$\text{rv}$
TTNNNN	1	is a function from the
TNTTNN		non-empty collection $C$ of
TNNNTN		subsets of $S'$ to the real
TNNNNT	2	number line $\mathbb{R}$ .
NTTTNN		In the T-S case study, the
NTNTTN		elements $s$ of $S'$ look like
NTNNNT		NNNTNN and the rv $\bar{Y}$
NNNTTN		counts how many Ts they contain.
NNNTNT		
NNNNTT		
:	:	
TTTTTT	5	

For instance,  $\Sigma(TNNTTN) = 2$  and (62)  
 $\Sigma(NNNNTT) = (0+0)/2$  (i.e.,  $\Sigma$  ignores  
the order of the children). We can

use the following notation to simplify things.

Notation   $P(\overset{\text{TF proportion}}{\Sigma} = \gamma) \stackrel{\leftarrow \text{set } \rightarrow}{=} P(\{s : \Sigma(s) = \gamma\})$

For example,  $P(\Sigma = 1) = P(\{s \in S : \Sigma(s) = 1\})$

$= P(\{TNNTTN, NTNNT, NN+NN, NNNTH,  
 NNNNT\}).$  In general the values  
 a random variable takes

can be just about anything, but  
 in this course all of our rvs will

be real-valued

In the T-S case study  
 the rv  $\Sigma$  can only take  
 on the values  $0, 1, \dots, 5.$

$y$	$P(I=y)$
0	0.237
1	0.396
2	0.264
3	0.088
4	0.015
5	0.001

You can see that a rv  $I$  (63) is completely specified by two things: the values it can take on, and the probability for those values.

(see p. (2))

Definition / The (probability)

distribution of a random variable

$I$  is the collection of all probabilities of the form  $P(I \in A)$  for all sets  $A$  of real numbers in the non-empty collection  $\mathcal{C}_R$  of subsets of the real

number line  $R$ . [The rv  $I$  in the

T-s we study has a finite set of possible values —

This is true of some, but not all, rvs. (64)

Definition A random variable has a discrete distribution, or equivalently  $\Sigma$  is a discrete rv, if the set of (distinct) values of  $\Sigma$  is finite or at most countably infinite; rvs for which the set of possible values is uncountable are called continuous random variables.

Example ① The rv  $\bar{X} = \begin{cases} 1 & \text{if } \Sigma > 0 \\ 0 & \text{otherwise} \end{cases}$  (with  $\Sigma = \# \text{Trials}$ ) is discrete, taking on only the values  $\{0, 1\}$  - such rvs are called dichotomous or binary.

② Imagine a scale for weighing things (65)  
that has a dial you can set to specify  
how many significant figures<sup>(sigfigs)</sup> of precision  
you want. Buy a "1 pound" package of  
butter at your favorite market and weigh it.

possible  
weights  
(ounces)

16  
16.0  
15.99  
15.993  
15.9928  
⋮

If there's no conceptual  
limit to the number of  
sigfigs you could get,  
a rv  $\Sigma = \begin{cases} \text{the actual (true)} \\ \text{weight of the package} \end{cases}$   
should be modeled as  
continuous, having values  
(e.g.) on  $(0, \infty)$ , the positive  
part of  $\mathbb{R}$ .

Reality check: Infinite  
precision is impossible in practice;

(66)

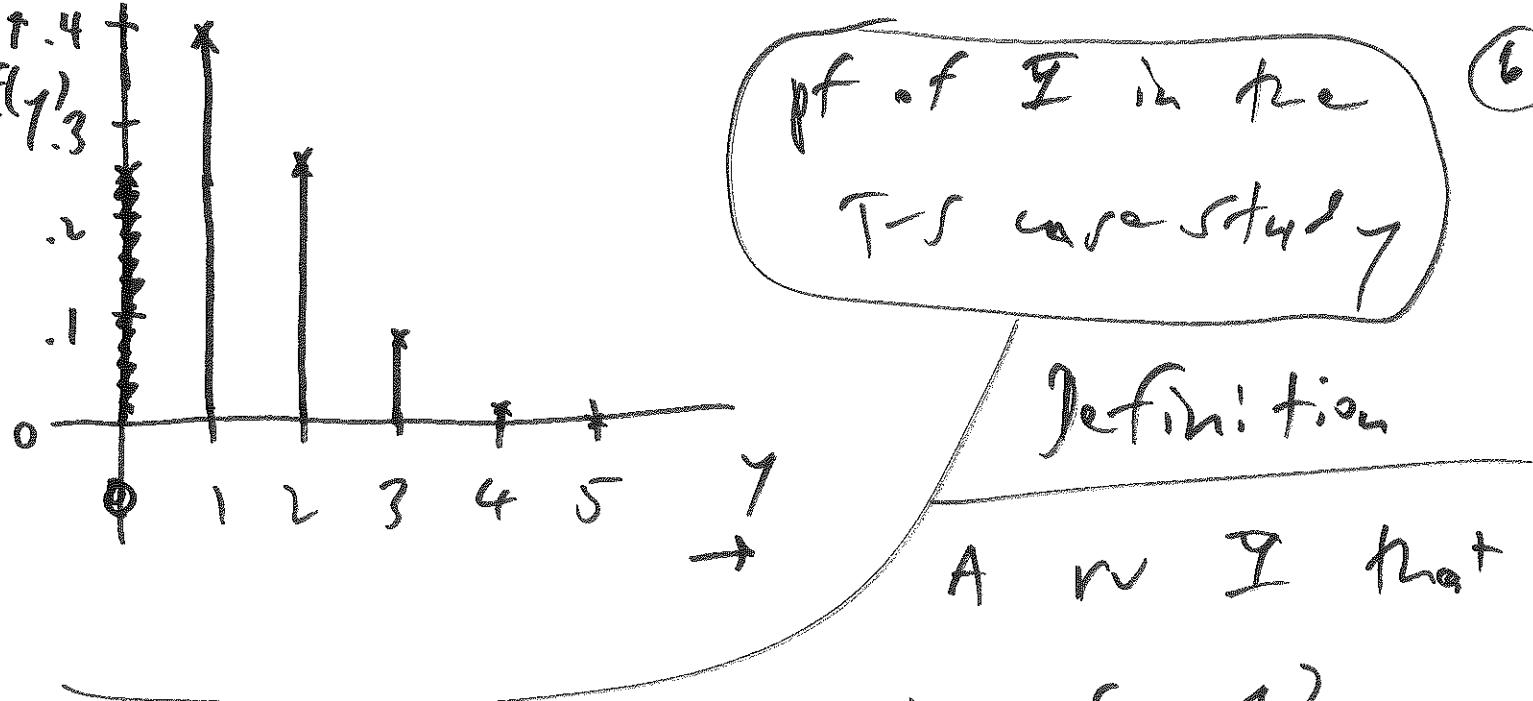
every measurement you ever make is in actuality discrete, but it's useful to regard rvs that are conceptually continuous (i.e., no limit in principle to the precision of measurement) as continuous.

Definition Given a

(rvs)

discrete rv  $\Sigma$ , the probability function (pmf or pf) of  $\Sigma$  is the function  $f$  that keeps track of the probabilities associated with  $\Sigma$ :  $f(y) = P(\Sigma = y)$ . The set  $\{y : f(y) > 0\}$  is called the support of (the distribution of)  $\Sigma$ .

(DS is almost unique in using pf, nearly everybody talks about the pmf.)



only takes on the values  $\{0, 1\}$  -

i.e., a binary rv - is said to have  
a Bernoulli distribution with

parameter  $p$  - written  $\text{Bernoulli}(p)$  -  
if 
$$f(y) = p(1-y) = \begin{cases} p & \text{for } y=1 \\ 1-p & y=0 \end{cases}$$

James  
Bernoulli  
Swiss (1655-  
1705)

Notation  $\left\{ \begin{array}{l} I \text{ follows} \\ \text{a Bernoulli}(p) \\ \text{distribution} \end{array} \right\} \leftrightarrow I \sim \text{Bernoulli}(p)$  is distributed as

Example] In the powerball lottery (see homework 1 problem 2) 5 white balls are drawn at random without replacement from a bin with balls numbered  $\{1, 2, \dots, 69\}$ . Let  $\underline{W}_i = \#$  on  $i^{\text{th}}$  drawn <sup>white</sup> ball.

$$\text{Clearly } p(\underline{W}_1 = w_1) = \begin{cases} \frac{1}{69} & w_1 = 1, 2, \dots, 69 \\ 0 & \text{otherwise} \end{cases}$$

less clearly (but true)  $\underline{W}_2, \dots, \underline{W}_5$  follow the same distribution if nothing is known about the previous draws.

Definition] For any two integers  $a \leq b$ , a  $w \in \Omega$  that's equally likely to be any of the values  $\{a, a+1, \dots, b\}$  has the uniform distribution Uniform  $\{a, b\}$ . Evidently,

$$\text{it is if } f(y) = P(\bar{\gamma} = y) = \begin{cases} \frac{1}{b-a+1}, & y = a, \dots, b \\ 0, & \text{else} \end{cases}$$

(69)

$$\bar{\gamma} \sim \text{Uniform } \{a, b\} \Leftrightarrow \bar{\gamma} \underset{\substack{\text{chosen at random} \\ \text{from } \{a, a+1, \dots, b\}}}{}$$

Definition  $\hookrightarrow$   $n$  trials  $\xrightarrow{\text{random}}$   $\xrightarrow{\text{(Aug 16)}}$

with each trial recorded as a success

S or failure F. If each trial is

$\xrightarrow{\text{not sample space}}$  independent of all the others and  
the chance  $P$  of success is constant

across the trials, then  $\bar{\gamma} = \# \cdot \text{successes}$

has the Binomial distribution

cptf

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

(with parameters  $n$  and  $p$ )

In shorthand  $\mathbf{I} \sim \text{Binomial}(n, p)$ . (70)

let  $B_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \dots \quad \text{failure} \end{cases}$

for  $i = 1, \dots, n$ ; then under these assumptions

$B_i \stackrel{\text{IID}}{\sim} \text{Bernoulli}(p)$  and all the  $B_i$  are

independent.

Notation  $\mathbf{X}_i \stackrel{\text{IID}}{\sim} f(x_i)$

means that all of the vars  $X_1, X_2, \dots$  are independent and identically distributed draws from the distribution with pf

$f(x_i)$ .

Thus with the success/failure

trials,  $B_i \stackrel{\text{IID}}{\sim} \text{Bernoulli}(p)$  and  
 $(i=1, \dots, n)$

$\mathbf{I} = \sum_{i=1}^n B_i \sim \text{Binomial}(n, p)$ .

This is our first example of the distribution of the sum of 4 bunches of IIS vs., a topic we'll examine in detail later.