

Extension of the LTP (4)

Assuming all conditional probabilities are defined

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

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

Definition

Events A, B are independent iff

← (freq)

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

which (as long as $P(A) > 0, P(B) > 0$)

is equivalent to

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

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

← (Bayesian)

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)
($j=2, \dots, k$),

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

(Bayesian)
Interpretation
of independence

A, B independent \iff

(54)

information about A

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

Definition

Another ^{useful} 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 $\{A_1, \dots, A_k\}$
($j = 2, \dots, k$)

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

← product

Statistical Example

Suppose that there is a machine that can ⁽⁵⁵⁾

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

Someone sets the dial to a θ value that's unknown to you and starts producing coin tosses I_1, I_2, \dots

Suppose the first 10 tosses come out
1 0 1 1 1 0 0 1 1 1 ← "bits" (binary digits)
HTHTHTTHTHH (7 H, 3 T) ↗
(John Tukey)

Q: Is there information in these first 10 tosses that helps you to predict I_n ?

A: Yes, definitely: it looks like (56)

θ is around $\frac{7}{10}$, so you would predict

$I_n = H$. Thus I_n depends on I_1, \dots, I_{n-1} probabilistically.

Now, suppose instead that you watched the ~~person~~ ^{person} with the machine

set the dial to $\theta = 0.81$, so that

θ is now known to you. The next 10

tosses come out H H H T H T H H H H

(8 H, 2 T). **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_n

that helps you to predict any of (57)

the other I_j .

given θ , the I_i are indep.

Thus the I_i are

unconditionally dependent but

conditionally independent given θ .

(4 Aug 17)

Bayes's Theorem for events
(a finite partition)

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)

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

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 .

(59)
Before ^(a priori) 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 ^(a posteriori) the dataset A
arrives, you can use Bayes's 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.

(the denominator)
 $P(A)$ 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. Thus

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

is interpreted as

$$\begin{matrix} \text{(posterior} \\ \text{information)} \end{matrix} = \begin{matrix} \text{(prior} \\ \text{information)} \end{matrix} \cdot \begin{matrix} \text{(data)} \\ \text{(likelihood} \\ \text{information)} \end{matrix} \\ \hline \begin{matrix} \text{(normalizing} \\ \text{constant)} \end{matrix} .$$

Random variables and their distributions (61)

Example: Tay-Sachs Disease

T = T-s baby
N = not

(S) 2

NNNNN	0
TNNNN	1
NTNNN	
NNTNN	
NNNTN	
NNNNT	
TTNNN	2
TNTNN	
TNNTN	
TNNNT	
NTTNN	
NTNTN	
NTNNT	
NNTTN	
NNTNT	
NNNTT	
⋮	⋮
TTTTT	5

← # of T-s babies = \underline{Y}

Given a sample

Definition

Space S for an experiment E ,
a (real-valued) random variable
(RV) is a function from the
non-void collection C of
subsets of S to the real
number line \mathbb{R} .

In the T-s case study, the
elements s of S look like
NNNTN and the RV Y
counts how many Ts they contain.

For instance, $\mathcal{I}(TNNTN) = 2$ and $\textcircled{62}$
 $\mathcal{I}(NNNTT) = (2/5) \cdot 2$ (i.e., \mathcal{I} ignores
the order of the children).

We can

use the following notation to simplify things.

Notation $P(\mathcal{I} = y) \stackrel{\text{IFF proposition}}{=} P(\{s \in \mathcal{S} : \mathcal{I}(s) = y\})$

For example, $P(\mathcal{I} = 1) = P(\{s \in \mathcal{S} : \mathcal{I}(s) = 1\})$
 $= P(\{TNNNN, NTNNN, NNTNN, NNNTN,$
 $NNNNT\})$.

In general the values
a random variable takes

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

be real-valued

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

(7 Aug 17)

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 of the (probability) distribution of a random variable

\mathcal{I} is the collection of all probabilities of the form $P(I \in A)$ for all sets A of real numbers in the σ -algebra collection $\mathcal{C}_{\mathbb{R}}$ of subsets of the real number line \mathbb{R} .

The rv I in the

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

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

Definition A random variable X has a discrete distribution, or equivalently

X is a discrete rv, if the set of (distinct) possible values of X 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 $X = \begin{cases} 1 & \text{if } Y > 0 \\ 0 & \text{otherwise} \end{cases}$

(with $Y = \# \text{ T-r balls}$) is discrete, taking on only the values $\{0, 1\}$ - such rvs are called dichotomous or binary.
{yes, no} {1, 0}

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

possible weights (pounds)

16

16.0

15.99

15.9930

15.9928

⋮

⋮

If there's no conceptual limit to the number of sigfigs you could get,

a rv $X =$ (the actual (true) weight of the package)

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;

every measurement you ever make is (66)
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.

(23 Apr 19) Definition Given a
discrete rv \mathcal{I} , the probability function
(pmf or pf) of \mathcal{I} is the function
 $f_{\mathcal{I}}$ that keeps track of the probabilities
associated with \mathcal{I} : $f_{\mathcal{I}}(y) = P(\mathcal{I} = y)$.
The set $\{y: f_{\mathcal{I}}(y) > 0\}$ is called the
support of (the distribution of) \mathcal{I} .

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