

4/23/19

Lecture 7

Law of total probability:

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

These probabilities
add up to 1

These don't

Expressing the probability of A as a weighted average of the conditional probabilities of A given the partition sets weighted by the probabilities of those partition sets.

Extension of the LTP: Assuming all conditional probabilities are defined in what follows, if C is

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

Def: Events A, B are independent iff

$$P(A \cap B) = P(A) \cdot P(B), \text{ which (as long as } P(A) > 0,$$

$P(B) > 0$) is equivalent to the Bayesian view that

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

Consequences of the def. of independence:

1) If A and B are independent, then so are A and B^c , A^c and B , and A^c and B^c

2) Extension of the def. to more than 2 events:

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

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

Bayesian Interpretation of Independence: A, B independent iff information about A doesn't change the chances associated with B , and vice versa

Def: 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 Ex:

Suppose that there is a machine that can take an ordinary coin and produce IID tosses of the coin with $P(H) = \theta$ with H being any one toss, 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 $\bar{Y}_1, \bar{Y}_2, \dots$

Suppose the first 10 tosses come out 1011100111

H T H H H T T H H H

(7 H, 3 T) "bits" (binary digits)

Q: Is there information in these first 10 tosses that helps you predict the next outcome \bar{Y}_{11} ?

A: Yes, definitely. It looks like θ is around $\frac{7}{10}$ so you would predict $\bar{Y}_{11} = H$. Thus \bar{Y}_{11} depends on \bar{Y}_1 through \bar{Y}_{10} probabilistically.

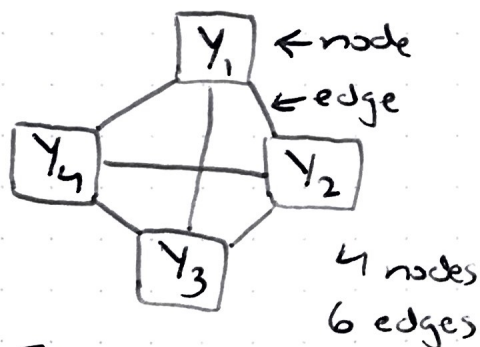
Now suppose instead that you watched the person with the machine set the dial to $\theta = 0.81$ so that θ is now known to you. The next 10 tosses come out HHTHTHTHTHH (8H, 2T)

Q: Is there information in these 10 tosses that helps you predict the next toss?

A: No. You know that $\theta = 0.81$ so there's no information in any of the \mathcal{Y}_i that helps you to predict any of the other \mathcal{Y}_j . Given θ , the \mathcal{Y}_i are independent.

Thus, the \mathcal{Y}_i are unconditionally dependent but conditionally dependent given θ .

Ex with 4 \mathcal{Y} values: Graph Theory



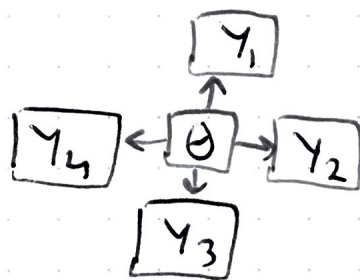
Edges denotes any dependence amongst the nodes

\mathcal{Y}_i dependent

n nodes \Leftarrow # of data points

$$\binom{n}{2} \text{ edges} = \frac{n(n-1)}{2} = O(n^2)$$

This is of order n^2
"Big O" of n^2



$x-y$	$x \rightarrow y$
x & y are dependent	x causes y

As soon as we know what θ is, the \mathcal{Y}_i become conditionally dependent given θ .

The past and future are conditionally independent given the truth (θ)

Bayes' Theorem for events based on the possibility of 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:

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

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 .

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

$P(B_i)$

(A posteriori) - 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 aka 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 Theorem to make all the $P(B_i | A)$ add up to 1

Thus $\frac{P(B_i)P(A|B_i)}{P(A)}$ is interpreted as

$$\begin{array}{c} \text{(posterior} \\ \text{information)} \end{array} = \frac{\begin{array}{c} \text{(prior} \\ \text{information)} \end{array} \cdot \begin{array}{c} \text{(data)} \\ \text{(likelihood} \\ \text{information)} \end{array}}{\begin{array}{c} \text{(normalizing} \\ \text{constant)} \end{array}}$$

Back to credit card case study:

We know $P(B) = 0.01$ $P(S-|G) = 0.97$ $P(S+|B) = 0.98$

We want $P(B|S+)$

Only 2 possible states: (B, G)

truth (unknown) \uparrow system says (data)

Method 1: 2x2 table

Method 2: Bayes' Theorem in odds form

$$\frac{P(B|S+)}{P(\text{not } B|S+)} = \left[\frac{P(B)}{P(\text{not } B)} \right] \cdot \left[\frac{P(S+|B)}{P(S+|\text{not } B)} \right]$$

We know all the probabilities on the right side

$$\begin{array}{c} \text{(posterior odds} \\ \text{in favor of} \\ B \end{array} = \begin{array}{c} \text{(prior odds} \\ \text{in favor} \\ \text{of } B \end{array} \cdot \begin{array}{c} \text{(Bayes factor} \\ \text{in favor of } B \end{array}$$

$$\left. \begin{array}{l} P(B) = 0.01 \\ P(\text{not } B) = P(G) = 0.99 \\ P(S+|B) = 0.98 \\ P(S+|\text{not } B) = 1 - P(S-|\text{not } B) = 1 - 0.97 = 0.03 \end{array} \right\} \text{So } \frac{P(B|S+)}{P(\text{not } B|S+)} = \left(\frac{0.01}{0.99} \right) \left(\frac{0.98}{0.03} \right)$$

$$= \left(\begin{array}{l} 99 \text{ to } 1 \text{ prior} \\ \text{odds against } B \end{array} \right) \left(\begin{array}{l} 98 \text{ to } 3 \text{ odds} \\ \text{in favor of } B \end{array} \right)$$

$$= \frac{98}{(99)(3)} = \frac{98}{297} \leftarrow \text{posterior odds in favor of } B$$

o = odds
p = probability

$$o = \frac{p}{1-p} \rightarrow p = \frac{o}{1+o}$$

$$p = \frac{\frac{98}{297}}{1 + \frac{98}{297}} = \frac{98}{98 + 297} = \frac{98}{395} = 0.25 \text{ as before}$$

Method 3: Bayes' Theorem in probability form

$$P(B|S+) = \frac{P(B) P(S+|B)}{P(S+)} = \frac{(0.01)(0.98)}{?}$$

Truth: {B, G} (unknown)

↑ prediction

Data: {S+, S-}

$P(B|S+)$ Search for truth on basis of data: statistical inference
 ↑ truth (unknown) ↑ data (observable)

unknown (true) state of world may well be unobservable

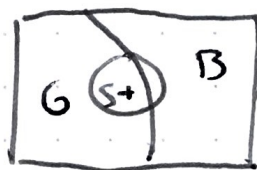
We don't know S+ by itself so we need B or G

When \mathcal{Y} is hard, get help by extending the conversation.

Find some other aspect of world \mathcal{X} upon which \mathcal{Y} depends and predict \mathcal{Y} in 2 stages: $\left\{ \begin{array}{l} \mathcal{X} \\ (\mathcal{Y}|\mathcal{X}) \end{array} \right\}$

$P(S+) = \text{data } (\mathcal{Y})$

Other info (\mathcal{X}) = {G, B}



Partition the world

$$P(S+) = P[(S+ \text{ and } G) \cup (S+ \text{ and } B)]$$

$$= P(S+ \text{ and } G) + P(S+ \text{ and } B)$$

There's no overlap

$$= P(G)P(S+|G) + P(B)P(S+|B)$$

$$= (0.99)(1 - P(S-|G)) + (0.01)(0.98)$$

$$= (0.99)(1 - 0.97) + (0.01)(0.98)$$

$$= (0.99)(0.03) + (0.01)(0.98)$$

$$= 0.0297 + 0.0098 = 0.0395$$

$$P(B|S+) = \frac{P(B)P(S+|B)}{P(S+)} = \frac{(0.01)(0.98)}{0.0395} = \frac{98}{395} = 0.25$$

Ch. 3 Random variables and their distributions

Ex: Tay-Sachs Disease

T = T-S baby

N = not T-S baby

Def: Given a sample space S

for an experiment E , a (real-valued) random variable (rv) is a function from the non-empty collection \mathcal{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 $NMNTN$ and the rv Y counts how many Ts they contain.

$S \rightarrow$	Y	# of T-S babies
NNNNN	0	= Y
TNNNN	1	
NTNNN		
NNTNN		
NNNTN		
NNNTN		
TTNNN		
TNTNN		
TNNTN		
TNNNT		
NNTTN	2	
NTNTN		
NTNNT		
NNTTN		
NNTNT		
NNNTT		
⋮	⋮	
TTTTT	5	

For instance, $Y(TNNTN) = 2$ and $Y(NNNTT) = 2$ also. Y ignores the order of the children.

To simplify: $P(Y=y) = P(\underbrace{\{s : Y(s) = y\}}_{\text{T/F proposition}})$

set

$$\text{For ex) } P(\bar{Y} = 1) = P(\{s \in S : \bar{Y}(s) = 1\})$$

$$= P(\{TNNNN, NTNNN, NNTNN, \\ NNNNT, NNNNT\})$$

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 \bar{Y} can only take on the values $0, 1, \dots, 5$

$$P(\bar{Y} = \gamma)$$

random variable (process) a possible value of \bar{Y} (outcome)

γ	$P(\bar{Y} = \gamma)$
0	0.237
1	0.396
2	0.264
3	0.088
4	0.015
5	0.001

You can see that a rv \bar{Y} is completely specified by 2 things: the values it can take on and the probabilities for those values.

Def: The probability distribution of a random variable \bar{Y} is in the collection of all probabilities of the form $P(\bar{Y} \in A)$ for all sets A of real #'s in the non-void collection $\mathcal{C}_{\mathbb{R}}$ of subsets of the real number line \mathbb{R}

The rv \bar{Y} in the T-S case study has a finite set of possible values — this time of same, but not all, rvs

Def: A random variable \bar{Y} has a discrete distribution, or equivalently \bar{Y} is a discrete rv if the set of possible distinct values of \bar{Y} is finite or at most countably infinite;

rvs for which the set of all possible values is uncountable are called continuous random variables

Ex: The rv $X = \begin{cases} 1 & \text{if } Y > 0 \\ 0 & \text{otherwise} \end{cases}$ (with $Y = \# \text{ T-S babies}$)

is discrete, taking on only the values $\{\text{yes, no}\} \rightarrow \{0, 1\}$
such rvs are called dichotomous or binary

Ex: Imagine a scale for weighing things that has a dial you can set to specify how many significant figures (sig figs) of precision you want.

Buy a "1 pound" package of butter and weigh it.

<u>Possible Weights</u>	If there's no conceptual limit to the number of sigfigs you could get, a rv $Y =$ (the actual true weight of the package) should be modeled as continuous, having positive values on the real number line $\mathbb{R} (0, \infty)$
16	
16.0	
15.99	
15.9930	
15.9928	
\vdots	

Reality Check: Infinite precision is impossible in practice, every measurement you ever makes 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