

Random Variables

Bayes Rule Applications

CSE 312 Spring 24
Lecture 9

Today

A convenient representation: random variables

Bayes' Rule in the real world!

Implicitly defining Ω

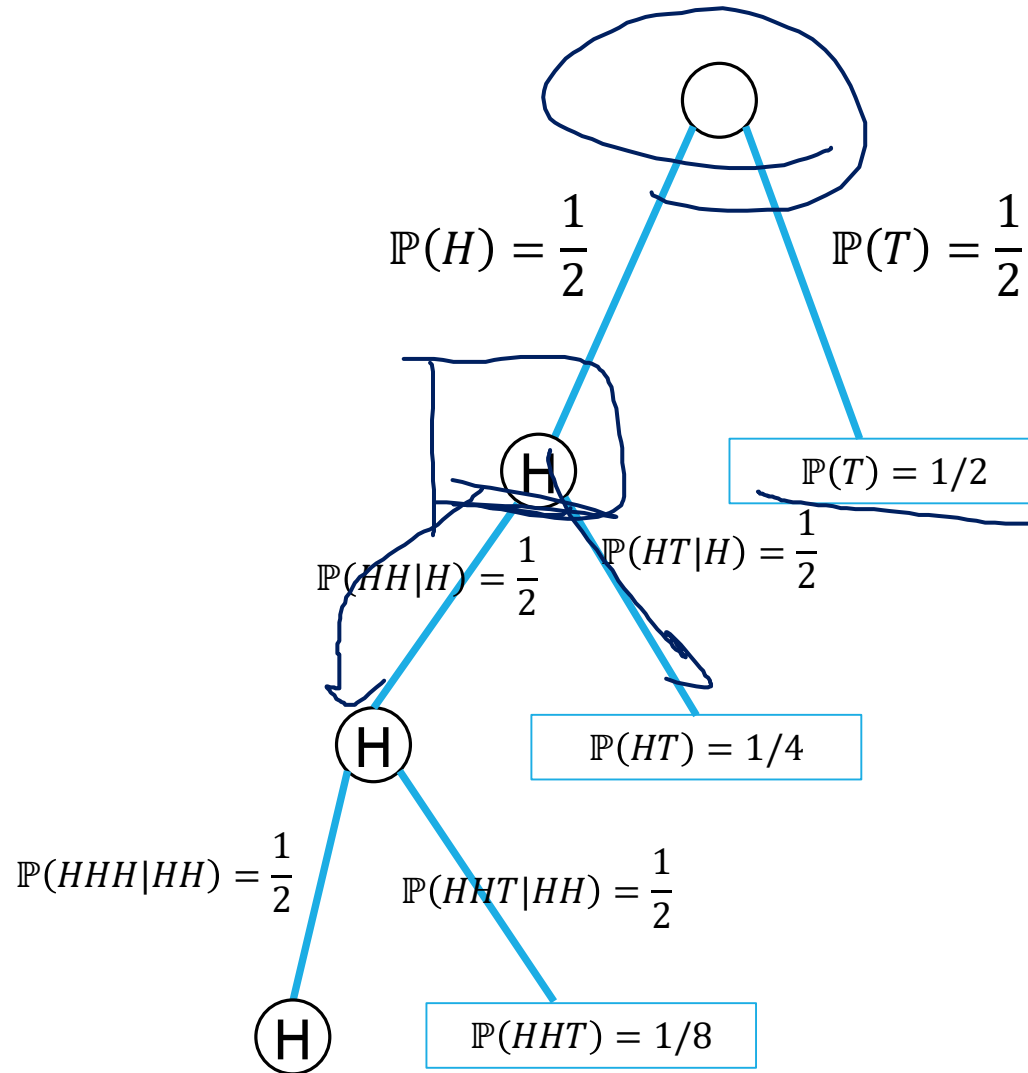
We've often skipped an explicit definition of Ω .

Often $|\Omega|$ is infinite, so we really couldn't write it out (even in principle).

How would that happen?

Flip a fair coin (independently each time) until you see your first tails.
what is the probability that you see at least 3 heads?

An infinite process.



Ω is infinite.

A sequential process is also going to be infinite...

But the tree is "self-similar"

To know what the next step looks like, you only need to look back a finite number of steps.

From every node, the children look identical (H with probability $\frac{1}{2}$, continue pattern; T to a leaf with probability $\frac{1}{2}$)

Finding $\mathbb{P}(\text{at least 3 heads})$

Method 1: infinite sum.

$\underbrace{HHHT}_{1/2^4}$
 $\underbrace{HHTT}_{1/2^5}$
 $\underbrace{HTTT}_{1/2^6}$

Ω includes $H^i T$ for every i . Every such outcome has probability $1/2^{i+1}$

What outcomes are in our event?

$$\sum_{i=3}^{\infty} 1/2^{i+1} = \frac{1/2^4}{1-1/2} = \frac{1}{8}$$

Infinite geometric series, where common ratio is between -1 and 1 has closed form $\frac{\text{first term}}{1-\text{ratio}}$

Finding \mathbb{P} (at least 3 heads)

Method 2:

Calculate the complement

$$\mathbb{P}(\text{at most 2 heads}) = \frac{1}{2} + \frac{1}{4} + \frac{1}{8}$$

$$\mathbb{P}(\text{at least 3 heads}) = 1 - \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} \right) = \frac{1}{8}$$





Random Variables

Random Variable

What's a random variable?

Formally

Random Variable

$X: \Omega \rightarrow \mathbb{R}$ is a random variable
 $X(\omega)$ is the summary of the outcome ω

Informally: A random variable is a way to **summarize** the important (numerical) information from your outcome.

The sum of two dice

EVENTS

We could define

E_2 = "sum is 2"

E_3 = "sum is 3"

...

E_{12} = "sum is 12"

And ask "which event occurs"?

RANDOM VARIABLE

$X: \Omega \rightarrow \mathbb{R}$

X is the sum of the two dice.

↪

More random variables

From one sample space, you can define many random variables.

Roll a fair red die and a fair blue die

Let D be the value of the red die minus the blue die $D(4,2) = 2$

Let S be the sum of the values of the dice $S(4,2) = 6$

Let M be the maximum of the values $M(4,2) = 4$

...

Notational Notes

$f(\text{input})$

We will always use capital letters for random variables.

It's common to use lower-case letters for the values they could take on.

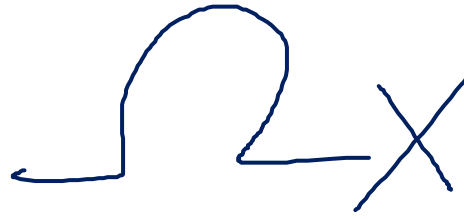
Formally random variables are functions, so you'd think we'd write

$$\underline{\underline{X(H, H, T) = 2}}$$

But we nearly never do. We just write $X = 2$

X

Support (Ω_X)



The "support" (aka "the range") is the set of values X can actually take.

We called this the "image" in 311.

D (difference of red and blue) has support $\{-5, -4, -3, \dots, 4, 5\}$

S (sum) has support $\{2, 3, \dots, 12\}$

What is the support of M (max of the two dice)

Probability Mass Function

Often we're interested in the event $\{\omega: X(\omega) = x\}$

$$P_S(7) = P(S=7) = \frac{1}{6}$$

Which is the event...that $X = x$.

$$X = 5$$

We'll write $\mathbb{P}(X = x)$ to describe the probability of that event

$$\text{So } \mathbb{P}(S = 2) = \frac{1}{36}, \mathbb{P}(S = 7) = \frac{1}{6}$$

The function that tells you $\mathbb{P}(X = x)$ is the "probability mass function"

We'll often write $p_X(x)$ for the pmf.

$$p_X(x)$$

Partition

A random variable partitions Ω .

Let T be the number of twos in rolling a (fair) red and blue die.

$$p_T(0) = 25/36$$

$$p_T(1) = 10/36$$

$$p_T(2) = 1/36$$

	D2=1	D2=2	D2=3	D2=4	D2=5	D2=6
D1=1	(1,1)	(1,2)	(1,3)	(1,4)	(1,5)	(1,6)
D1=2	(2,1)	(2,2)	(2,3)	(2,4)	(2,5)	(2,6)
D1=3	(3,1)	(3,2)	(3,3)	(3,4)	(3,5)	(3,6)
D1=4	(4,1)	(4,2)	(4,3)	(4,4)	(4,5)	(4,6)
D1=5	(5,1)	(5,2)	(5,3)	(5,4)	(5,5)	(5,6)
D1=6	(6,1)	(6,2)	(6,3)	(6,4)	(6,5)	(6,6)

Try It Yourself

$$P_X(x) = P(X=x)$$

There are 20 balls, numbered 1,2,...,20 in an urn.

You'll draw out a size-three subset. (i.e. without replacement)

$\Omega = \{\text{size three subsets of } \{1, \dots, 20\}\}$, $\mathbb{P}()$ is uniform measure.

Let X be the largest value among the three balls.

If outcome is $\{4,2,10\}$ then $X = 10$.

Write down the pmf of X

$$P_X(7) = P(X=7)$$

$$\frac{\binom{6}{2}}{\binom{20}{3}}$$

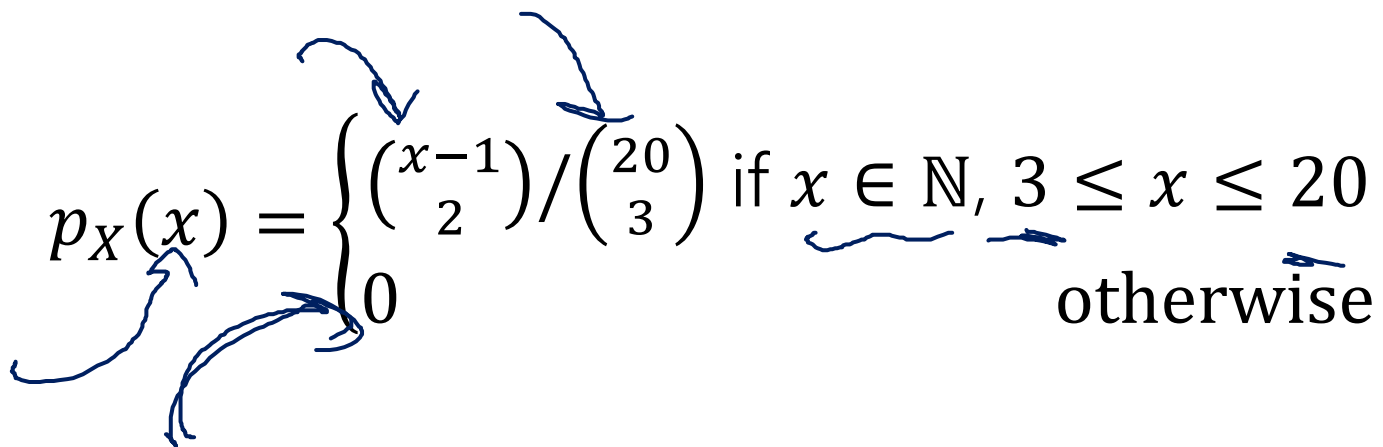
Fill out the poll everywhere so Robbie knows how long to explain
Go to pollev.com/robbie

Try It Yourself

There are 20 balls, numbered 1,2,...,20 in an urn.

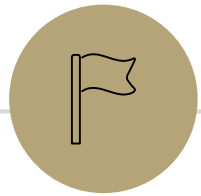
You'll draw out a size-three subset. (i.e. without replacement)

Let X be the largest value among the three balls.

$$p_X(x) = \begin{cases} \binom{x-1}{2} / \binom{20}{3} & \text{if } x \in \mathbb{N}, 3 \leq x \leq 20 \\ 0 & \text{otherwise} \end{cases}$$
Handwritten blue arrows point to the function $p_X(x)$, the binomial coefficient $\binom{x-1}{2}$, the binomial coefficient $\binom{20}{3}$, the condition $x \in \mathbb{N}, 3 \leq x \leq 20$, the value 0, and the word "otherwise".

Good check: if you sum up $p_X(x)$ do you get 1? 

Good check: is $p_X(x) \geq 0$ for all x ? Is it defined for all x ? 



Bayes in the real world



Application 1: Medical Tests

Helping Doctors and Patients Make Sense of Health Statistics

A researcher posed the following scenario to a group of 160 doctors:

Assume you conduct a disease screening using a standard test in a certain region. You know the following information about the people in this region:

The probability that a person has the disease is 1% (prevalence)

If a person has the disease, the probability that she tests positive is 90% (sensitivity)

If a person does not have the disease, the probability that she nevertheless tests positive is 9% (false-positive rate)

A person tests positive. She wants to know from you whether that means that she has the disease for sure, or what the chances are. What is the best answer?

A. The probability that she has the disease is about 81%.

B. Out of 10 people with a positive test, about 9 have the disease.

C. Out of 10 people with a positive test, about 1 have the disease.

D. The probability that she has the disease is about 1%

Let's do the calculation!

Let D be "the patient has the disease", T be the test was positive.

$$\mathbb{P}(D|T) = \mathbb{P}(T|D) \cdot \mathbb{P}(D) / \mathbb{P}(T)$$

$$= \frac{.9 \cdot .01}{.99 \cdot .09 + .01 \cdot .9} \approx 0.092$$

Calculation tip: for Bayes' Rule, you should see one of the terms on the bottom exactly match your numerator (if you're using the LTP to calculate the probability on the bottom)

Pause for vocabulary

Physicians have words for just about everything

Let D be has the disease; T be test is positive

$\mathbb{P}(D)$ is "prevalence"

$\mathbb{P}(T|D)$ is "sensitivity"

A 'sensitive' test is one which picks up on the disease when it's there (high sensitivity -> few false negatives)

$\mathbb{P}(\bar{T}|\bar{D})$ is "specificity"

A 'specific' test is one that is positive specifically because of the disease, and for no other reason (high specificity -> few false positives)

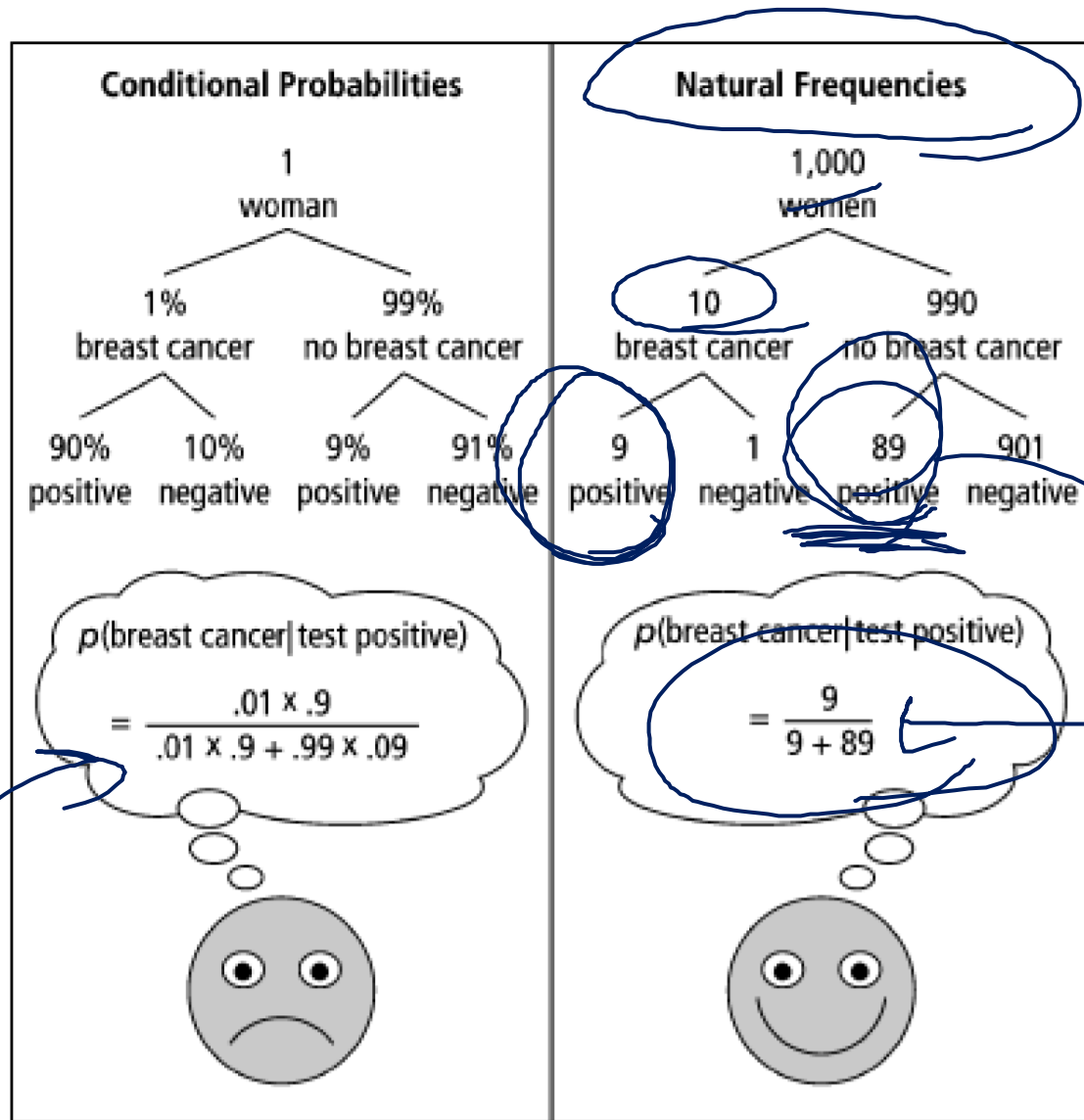
How did the doctors do

C (about 1 in 10) was the correct answer.

Of the doctors surveyed, less than $\frac{1}{4}$ got it right (so worse than random guessing).

After the researcher taught them his calculation trick, more than 80% got it right.

One Weird Trick!



Calculation Trick: imagine you have a large population (not one person) and ask how many there are of false/true positives/negatives.

What about the real world?

When you're older and have to do more routine medical tests, don't get concerned (yet) when they ask to run another test.*

It's usually fine.*

*This is not medical advice, Robbie is not a physician.



Careful Surveys



Application 2: An Imbalanced Survey

In 2014, a paper was published

[“Do non-citizens vote in U.S. elections?”](#)

This is a real paper (peer-reviewed). It claims that

1. In a survey, about 4% (of a few hundred) of non-U.S.-citizens surveyed said they voted in the 2008 federal election (which isn't allowed).
2. Those non-citizen voters voted heavily (estimate 80+%) for democrats.
3. “It is likely though by no means certain that John McCain would have won North Carolina were it not for the votes for Obama cast by non-citizens”

Application 2: What is this survey?

The “Cooperative Congressional Election Study” was run in 2008 and 2010.

It interviews about 20,000 people about how/whether they voted in federal elections.

Two strange observations:

1. The noncitizens are a very small portion of those surveyed. Feels a little strange.
2. Those people...maybe accidentally admitted to a crime?

Application 2: Another Red Flag

A response paper (by different authors)

[“The perils of cherry picking low frequency events in large sample surveys”](#)

Table 1

Response to citizenship question across two-waves of CCES panel.

Response in 2010	Response in 2012	Number of respondents	Percentage
Citizen	Citizen	18,737	99.25
Citizen	Non-Citizen	20	0.11
Non-Citizen	Citizen	36	0.19
Non-Citizen	Non-Citizen	85	0.45

An Explanation

Suppose 0.1% of people check the wrong check-box on any individual question (independently)

Suppose you really interviewed 20,000 people, of whom 300 are really non-citizens (none of whom voted), and the rest are citizens, of whom 70% voted. What is the probability someone appears to have voted

$$\mathbb{P}(\text{say } V | \text{say } NC) = \frac{\mathbb{P}(\text{say } NC | \text{say } V) \cdot \mathbb{P}(\text{say } V)}{\mathbb{P}(\text{say } NC)} = \frac{.001 \cdot .7}{.999 \cdot \left(\frac{300}{20000}\right) + .001 \cdot \left(\frac{19700}{20000}\right)} \approx 4.38\%$$

Conclusion

The authors of the original paper did know about response error...

...and they have an appendix that argues the population of “non-citizen” voters isn’t distributed exactly like you’d expect.

But with it being such a small number of people, this isn’t surprising.

And even they admit response bias played more of a role than they initially thought.

Though they still think they found some evidence of non-citizens voting (but not enough to flip North Carolina anymore).

Takeaways

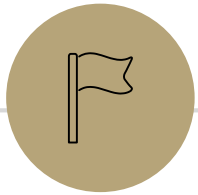
When talking about rare events (rare diseases, rare prize-winning-golden-tickets), think carefully about whether a test is really as informative as you think.

Do the explicit calculation

Intuition is easier if thinking about a large population of repeated tests, not just one.

Be careful of small subparts of large datasets

People from a large majority group (accidentally) clicking the wrong demographic information can “drown out” signal of a very small group.



Optional: Bayes Factor

A way to estimate Bayes calculations quickly



Bayes Factor

Another Intuition Trick: [from 3Blue1Brown](#)

When you test positive, you (**approximately**) multiply the prior by the “Bayes Factor” (aka likelihood ratio)

$$\frac{\text{sensitivity}}{\text{false positive rate}} = \frac{1 - FNR}{FPR}$$

Bayes Factor

Does it work?

Let's try it...

Find

$$\text{prior} \cdot \frac{\text{Sensitivity}}{FPR}$$

Wonka Bars

Willy Wonka has placed golden tickets on 0.1% of his Wonka Bars.

You want to get a golden ticket. You could buy a 1000-or-so of the bars until you find one, but that's expensive...you've got a better idea!

You have a test – a very precise scale you've bought.

If the bar you weigh **does** have a golden ticket, the scale will alert you 99.9% of the time.

If the bar you weigh does not have a golden ticket, the scale will (falsely) alert you only 1% of the time.

If you pick up a bar and it alerts, what is the probability you have a golden ticket?

Wonka Bars

Bayes Factor

$$\frac{99.9}{1}$$

Prior: .1%

Product: 9.99, so about 10%

About what Bayes Rule gets!

Application 1: Medical Tests

Helping Doctors and Patients Make Sense of Health Statistics

A researcher posed the following scenario to a group of 160 doctors:

Assume you conduct a disease screening using a standard test in a certain region. You know the following information about the people in this region:

The probability that a person has the disease is 1% (prevalence)

If a person has the disease, the probability that she tests positive is 90% (sensitivity)

If a person does not have the disease, the probability that she nevertheless tests positive is 9% (false-positive rate)

A person tests positive. She wants to know from you whether that means that she has the disease for sure, or what the chances are. What is the best answer?

A. The probability that she has the disease is about 81%.

B. Out of 10 people with a positive test, about 9 have the disease.

C. Out of 10 people with a positive test, about 1 have the disease.

D. The probability that she has the disease is about 1%

Bayes Factor

What about with the doctors?

$$1\% \cdot \frac{90\%}{9\%} = 10\%$$

Again about right!

Caution

Multiplying by the Bayes Factor is an **approximation**

It gives you the exact numerator for Bayes, but the denominator is "the number of false positives if the prevalence (/prior) were 0"

When the prior is close to 0, this is a fine approximation!

But plug in a prior of 15% on the last slide, and we get 150% chance.

What about negative tests?

For negative tests, the Bayes Factor is

$$\frac{FNR}{\text{specificity}}$$

Specificity is (1 – false negative rate)