

STAT 88: Lecture 4

Contents

Section 2.3: Bayes' Rule

Section 2.4: Use and Interpretation

Warm up: Let B_i be the event that a black card appears at Position i and R_i be the event that a red card appears at Position i .

(a) If you deal 2 cards, what is the chance the 2nd card is red? i.e. find $P(R_2)$.

(b) Find $P(R_{20} \cap R_{33})$, $P(R_{20} \cap B_{33})$, $P(B_{52} | R_{21} R_{40})$.

20th card is red

33rd card is black

$$(a) P(R_2) = P(R_1) = \frac{26}{52} \quad (\text{by symmetry})$$

$$(b) P(R_{20} \cap R_{33}) = P(R_{20}) P(R_{33} | R_{20}) \\ = \frac{26}{52} * \frac{25}{51}$$

$$P(R_{20} \cap B_{33}) = P(R_{20}) P(B_{33} | R_{20}) \\ = \frac{26}{52} * \frac{26}{51}$$

$$P(B_{52} | R_{21} \cap R_{40}) = \frac{26}{50}$$

Last time

Sec 2.1 The Chance of an Intersection

Multiplication (AND) rule:

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

Inclusion Exclusion (OR) rule:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

Sec 2.2 Symmetry in Simple Random Sampling

When randomly sampling from a population (with or without replacement), be aware of the difference between unconditional and conditional probability.

Example: Consider a deck of cards.

- $P(B_2) = 26/52$.
- $P(B_2|R_1) = 26/51$.

2.3. Bayes' Rule

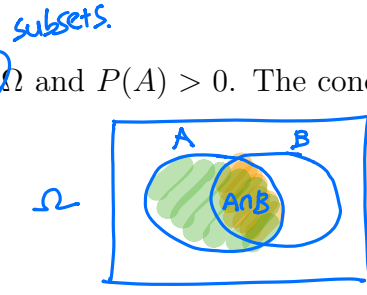
So far we have used the multiplication rule $P(A \cap B) = P(B|A)P(A)$ only in the settings where we can calculate $P(B|A)$ directly, e.g. $P(R_2|B_1) = 26/51$. Here we have 2 stages and we condition on what happens in the first stage.

The general definition of conditional probability, regardless of setting, is just a rearrangement of the multiplication rule.

Conditional Probability (Division Rule) Let $A, B \subseteq \Omega$ and $P(A) > 0$. The conditional probability of B given A is defined as

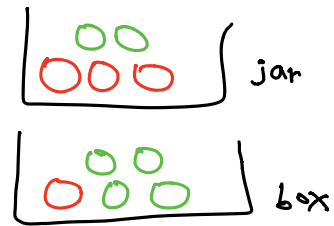
$P(B|A)$
↑
"given"

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$



Example: (a random container) I have two containers: a jar and a box. Each container has five balls.

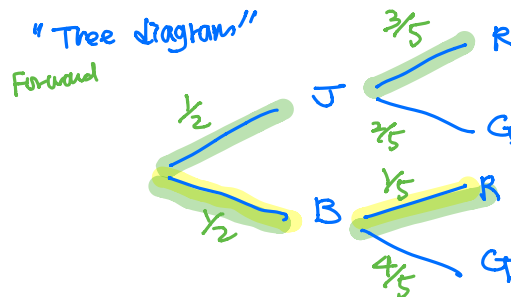
- The jar contains three red balls and two green balls.
- The box contains one red ball and four green balls.



1st stage 2nd stage.
Randomly pick a container and then a ball from that container. Given that the ball is red, what is the chance you pick the box?

Let J = jar
 B = box
 R = red
 G = green

Find $P(B|R)$
backward



$$\begin{aligned} P(J) &= \frac{1}{2}, & P(B) &= \frac{1}{2} \\ P(R|J) &= \frac{3}{5}, & P(G|J) &= \frac{2}{5} \\ P(R|B) &= \frac{1}{5}, & P(G|B) &= \frac{4}{5} \end{aligned}$$

$$\begin{aligned} P(B|R) &= \frac{P(B \cap R)}{P(R)} = \frac{P(B \cap R)}{P(B \cap R) + P(J \cap R)} \\ &= \frac{P(B) P(R|B)}{P(B) P(R|B) + P(J) P(R|J)} \\ &= \frac{\frac{1}{2} * \frac{1}{5}}{\frac{1}{2} * \frac{1}{5} + \frac{1}{2} * \frac{3}{5}} = \frac{1}{4} \end{aligned}$$

$$\begin{aligned} P(A \cap B) &= P(A|B) P(B) \\ &= P(B|A) P(A) \end{aligned}$$

posterior.
"updated prior given data"

We have updated our opinion about whether you picked the box or jar.

- Before we knew the color of the ball, we said the chance of drawing the box is $P(B) = 0.5$. This is called the **prior probability** of drawing the box.
- After we saw that the ball is red, we said that the chance that the box was drawn is 0.25. This is called the **posterior probability** of drawing the box.

$P(B|R)$

This way of updating probabilities based on new information is the basis for much inference in data science.

Bayes' Rule For $A, B \subseteq \Omega$,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B|A)}{P(A)P(B|A) + P(A^c)P(B|A^c)}.$$

Note that this is just the division rule in a particular scenario.

The rule helps you find a posterior probability: a conditional chance for the first stage, given the result of a second stage. The calculation has two ingredients:

- Probabilities for the first stage; these are called **prior probabilities**.
- Conditional probabilities for the second stage given the first; these are called **likelihoods**.

" $P(B|A)$

$P(A|B)$: posterior probabilities

Example: (Exercise 2.6.9) A factory has two widget-producing machines. Machine I produces 80% of the factory's widgets and Machine II produces the rest. Of the widgets produced by Machine I, 95% are of acceptable quality. Machine II is less reliable – only 85% of its widgets are acceptable.

Suppose you pick a widget at random from those produced at the factory.

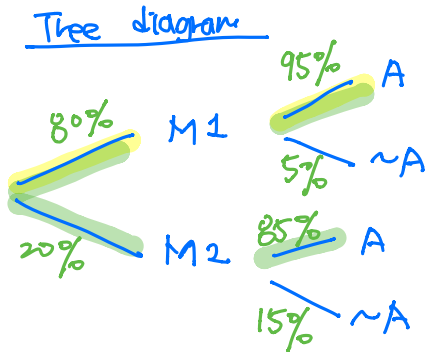
A = event of "acceptable". $M1$ = Machine 1. $M2$ = Machine 2.

1. Find the chance that the widget is acceptable, given that it is produced by Machine I.

$$P(A|M1) = 95\%$$

2. Find the chance that the widget is produced by Machine I, given that it is acceptable.

$$P(M1|A) = \frac{P(M1)P(A|M1)}{P(M1)P(A|M1) + P(M2)P(A|M2)}$$



$$= \frac{0.8 * 0.95}{0.8 * 0.95 + 0.2 * 0.85}$$

$$= \dots$$

2.4. Use and Interpretation

Harvard Medical School Survey (60 participants):

"If a test to detect a disease whose prevalence is 1/1,000 has a false positive rate of 5 per cent, what is the chance that a person found to have a positive result actually has the disease, assuming that you know nothing about the person's symptoms or signs?"

Harvard medical student answers ranged from 2% to 95%, with 27 out of the 60 Medical School members surveyed answering 95%. What do you say?

(Assume true pos rate is 1)

Terminology:

D = disease. + = test positive, - = test negative.

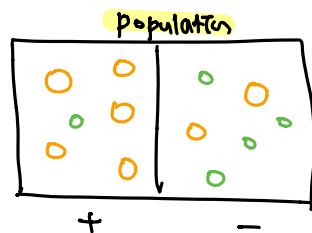
- **Prevalence** (also called the **base rate** of the disease) in the population is the percent of people who have the disease. $P(D) = \frac{1}{1000}$

- **True positive rate** is the rate of positive results among those who do have the disease. $P(+|D) = 1$

- **False positive rate** is the proportion of positive results among people who don't have the disease. $P(+|\sim D) = 0.05$

positive result - according to test the person has the disease

negative result - according to test the person doesn't have the disease.



○ healthy
○ sick

$$TPR = P(+|D) = \frac{10}{10}$$

$$FPR = P(+|\sim D) = \frac{1}{5}$$

D = disease

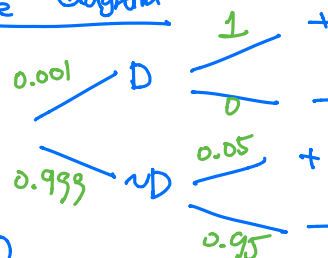
$\sim D$ = not disease (healthy)

+ = test positive

- = test negative.

Find $P(D|+)$

Tree diagram



$$P(D|+) = \frac{P(D) * P(+|D)}{P(D) * P(+|D) + P(\sim D) * P(+|\sim D)}$$

$$= \frac{0.001 * 1}{0.001 * 1 + 0.999 * 0.05} \approx 0.02$$

2%

Is this surprising?

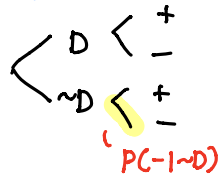
Yes and no.

Yes if $0.02 = 2\%$ seems small to you, and you were expecting a bigger number since true-positive rate is 1.

No. $P(D) = \frac{1}{1000}$. Suppose population consists of 100,000 people.

	+	-
D	100	0
$\sim D$	4995	94905

Base Rate Fallacy



$$P(D|+) = \frac{P(+|D)P(D)}{P(+)}$$

How did so many people get 95%?

$$P(-|\sim D) = 95\%$$

people confuse $P(D|+)$ for $P(-|\sim D)$

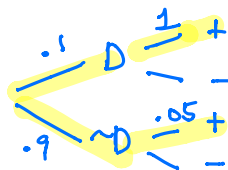
To compute $P(D|+)$ you need to take into account the base rate and people tend not to do this, instead focusing on likelihoods such as $P(-|\sim D)$

key point Posterior probability is effectively by base rate as well as the likelihoods

Suppose you have a 10% chance of having the disease because you show some symptoms or have a family history.

This changes prevalence to .1 from .001.

Now



$$P(D|+) = \frac{0.1 * 1}{0.1 * 1 + 0.9 * 0.05} = 0.69$$

Example: A True/False test consists of 60 questions. A student knows the answers to 45 of the questions. The remaining 15 answers he guesses at random by tossing a fair coin each time. If it lands heads he answers True and if it lands tails he answers False.

A question is picked at random from the 60 questions on the test. Given that the student got the right answer, what is the chance that he knew the answer?