

Unit 2– Introduction to Probability
Practice Problems

SOLUTIONS

1. Let A and B denote two independent genetic traits. Suppose the probability that an individual will exhibit trait A is $\frac{1}{2}$ and the probability that an individual will exhibit trait B is $\frac{3}{4}$. What is the probability that an individual will exhibit

(a) Both traits?

Answer: .375

$$\Pr[\text{both traits}] = \Pr[A]\Pr[B] = [.50][.75] = .375$$

(b) Neither trait?

Answer: .125

$$\Pr[\text{neither trait}] = \Pr[\text{not A}]\Pr[\text{not B}] = [.50][.25] = .125$$

(c) trait A but not trait B?

Answer: .125

$$\Pr[A \text{ and not B}] = \Pr[A]\Pr[\text{not B}] = [.50][.25] = .125$$

(d) trait B but not trait A?

Answer: .375

$$\Pr[\text{not A and B}] = \Pr[\text{not A}]\Pr[B] = [.50][.75] = .375$$

(e) exactly one trait?

Answer: .50

We sum the probabilities of the two mutually exclusive ways that yield “exactly one”

$$\begin{aligned} \Pr[\text{exactly one}] &= \Pr[(A, \text{not B}) \text{ or } (\text{not A}, B)] \\ &= \Pr[A, \text{not B}] + \Pr[\text{not A}, B] \\ &= [(.50)(.25)] + [(.50)(.75)] = .125 + .375 \\ &= .50 \end{aligned}$$

2. Suppose you are told that $\text{pr}(\text{right eye is blue}) = 1/3$ and $\text{pr}(\text{left eye is blue}) = 1/3$. Using the concepts and formulae in the lecture notes for Unit 2 (Introduction to Probability), confirm for yourself what you know by intuition, namely that $\text{pr}(\text{person is blue eyed}) = 1/3$ by solving for $\text{pr}(\text{blue right eye and blue left eye})$.

Under the assumption that left and right eye colors are always the same,

$\text{Pr}[\text{left is blue AND right is blue}]$ is the same as $\text{Pr}[\text{left is blue}] = 1/3$

$$1/3 = \text{Pr}[\text{left is blue and right is blue}]$$

$$= \text{Pr}[\text{right is blue}] \text{Pr}[\text{left is blue} | \text{right is blue}]$$

$$= \text{Pr}[\text{right is blue}] \{ 1 \}$$

$$= \text{Pr}[\text{right is blue}] \quad \checkmark$$

3. A physician develops a diagnostic test that is positive for 95% of the patients who have disease and is positive for 10% of the patients who do not have disease. Of patients tested, 20% actually have disease. Suppose you evaluate a patient by administering this diagnostic test and obtain a positive result. Using the information given, calculate the probability that this patient has disease.

Answer: .7037

Solution:

We want to calculate Probability (Disease | + test)

- **Probability (+ test | disease) = .95**

- **Probability (+ test | no disease) = .10**

Probability (Disease) = .20

Probability (not Disease) = .80

$$\begin{aligned} \text{Pr}(\text{disease} | +) &= \frac{\text{Pr}(\text{disease and } +)}{\text{Pr}(+)} && \text{by definition of conditional Probability} \\ &= \frac{\text{Pr}(+ | \text{disease}) \text{Pr}(\text{disease})}{\text{Pr}(+)} && \text{because we can re-write the numerator this way} \\ &= \frac{\text{Pr}(+ | \text{disease}) \text{Pr}(\text{disease})}{\text{Pr}(+ | \text{disease}) \text{Pr}(\text{disease}) + \text{Pr}(+ | \text{no disease}) \text{Pr}(\text{no disease})} \\ &= \frac{(.95) (.20)}{(.95) (.20) + (.10) (.80)} && = .7037 \end{aligned}$$

4. In introductory epidemiology, one of the study designs that are introduced is the **prospective cohort study**. In this type of study involving two groups, the investigator enrolls pre-set and known numbers of participants into each of the two groups that are generically described as “exposed” and “not exposed” and follows them forward to a designated end of the observation period, at which point some outcome is measured.

Consider the following prospective cohort study. A total of 1500 *never smoker* consenting heart attack survivors aged 60-65 are enrolled as “non-exposed”. An equal number, 1500 *current smoker* heart attack survivors aged 60-65 are enrolled as “exposed”. All are followed for a full 10 years and the occurrence of death recorded. Following are the data.

		Vital Status at 10 Years		
		Dead	Alive	
Exposure	Current Smoker	40	1460	1500
Status	Never Smoker	10	1490	1500
		50	2950	3000

- (a) Is it possible to estimate the probability of 10 year survival on the basis of these data?
Answer: Yes, but the question is a poor one as it does not specify time zero nor among whom. $2950/3000=0.9833$
- (b) Is it possible to estimate the relative risk of 10 year mortality that is associated with current cigarette use?
Answer: Yes, but this question too is poor. Without a meaningful time zero, the interpretation of the answer is non-existent. $(40/1500)/(10/1500)=4$
- (c) Is it possible to estimate the probability that a randomly selected person with vital status of “Alive” at 10 years is a current smoker?
Answer: Yes. Even though the study design called for fixed numbers of enrollments of current smokers and never smokers, it is possible to estimate the conditional probability of “current smoker” for a randomly selected person from among the vital status at 10 years of “Alive”.
- (d) Using these data, estimate the relative risk of 10 year mortality that is associated with current cigarette use.
Answer: $(40/1500)/(10/1500)=4$
- (e) Using these data, estimate the relative odds of 10 year mortality that is associated with current cigarette use. (Note – This question is asking you to compute an *odds ratio*).
**Answer: The event of interest is 10 year mortality. For this event,
 $OR = (Odds\ among\ current\ smokers)/(Odds\ among\ never\ smokers)$
 $= (40/1460)/(10/1490)=4.082$**

(f) Using these data, estimate the relative odds of a *current smoker* notation for non-survivors relative to survivors.

**Answer: Here, the event of interest is current smoker. For this event,
 $OR = (Odds\ among\ non-survivors)/(Odds\ among\ survivors)$
 $= (40/10)/(1460/1490)=4.082$**

(g) What do you notice about your answers to “e” and “f”?

**Answer: They are the same!
 This confirms what is noted in the lecture notes, that the OR for the event of disease in a cohort study comparison of exposure groups is equal to the OR for the event of history of exposure in a case-control comparison of disease/non disease groups.**

(h) How do your answers to “e” and “f” compare to your answer to “d”?

Answer: The OR is bigger than the RR; for rare diseases, the discrepancy becomes negligible..

5. Another study design that is introduced in introductory epidemiology is the **retrospective case-control study**. This is, by definition, a study that compares two groups. Here, the investigator enrolls pre-set and known numbers of participants into each of the two groups defined by disease status; “cases” are the enrollees with disease, “controls” do not have the disease under investigation. Retrospective review of the histories of all study participants is performed to identify the subsets in each of the case and control groups who have a history of the exposure of interest.

Consider the following retrospective case-control study of the association between coffee consumption and tumors of the lower urinary tract. The investigator enrolls 30 consenting cases that are patients with one or more tumors of the lower urinary tract. For comparison purposes, he/she also enrolls 100 consenting controls who have no such tumors. Following are the data.

		Tumors of Lower Urinary Tract		
		Yes	No	
History of Coffee consumption	5+ cups/day	20	44	64
	<1 cup/day	10	56	66
		30	100	130

(a) Is it possible to estimate the probability of one or more tumors of the lower urinary tract on the basis of these data?

Answer: No, because the case-control study design calls for fixed numbers of enrollment into the “tumors=yes” and “tumors=no” groups.

(b) Is it possible to estimate the relative risk of one or more tumors of the lower urinary tract that is associated with consumption of 5 or more cups of coffee/day?

Answer: No, for a reason that relates to the answer to “a”. It is not possible to estimate the component probabilities.

- (c) Using these data, estimate the relative odds of **high coffee consumption (5+ cups/day)** among cases, relative to controls.

Answer: For the event high coffee consumption,

$$\text{OR} = (\text{Odds among "tumor=yes"}) / (\text{Odds among "tumors=no"}) \\ = (20/10) / (44/56) = 2.545$$

- (d) Using these data, estimate the relative odds of **tumors of the lower urinary tract** among high coffee consumers (5+ cups/day), relative to non-coffee drinkers.

Answer: For the event "tumors=yes",

$$\text{OR} = (\text{Odds among "high coffee"}) / (\text{Odds among "non-coffee"}) \\ = (20/44) / (10/56) = 2.545$$

- (e) What do you notice about your answers to "c" and "d"?

Answer: They are the same.

6. Now consider a fully **cross-sectional study design**, this time with generic counts "a", "b", "c", and "d". In this design, the investigator does not do any formal enrollment. Counts are accumulated by observation. CDC surveillance programs are examples

		Disease	
		Yes	No
History of Exposure	Yes	a	b
	No	c	d

- (a) Using the letters "a", "b", "c", and "d", what is the formula for estimating relative odds of the **event of exposure** for persons with disease, compared to that for persons without disease?

Answer: [a/c] / [b/d] = (ad)/(bc).

- (b) Using the letters "a", "b", "c", and "d", what is the formula for estimating relative odds of the **event of disease** for exposed persons, compared to that for non-exposed persons?

Answer: [a/b] / [c/d] = (ad)/(bc).

- (c) Using the letters "a", "b", "c", and "d", what is the formula for estimating relative risk of the **event of disease** for exposed persons, compared to that for non-exposed persons?

Answer: [a/(a+b)] / [c/(c+d)].

- (d) What happens to your formula in your answer to #3c when the counts of disease (a and c) are very very small? Comment.

Answer: As "a" gets smaller and smaller, (a+b) → b

As "c" gets smaller and smaller, (c+d) → d. As a result,

$$\text{RR} = [a/(a+b)] / [c/(c+d)] \rightarrow [a/b] / [c/d] = (ad)/(bc) = \text{OR}.$$

So ... it is sometimes possible to estimate RR from a case-control study because

OR_{event=disease} = OR_{event=exposure}

When the disease is rare, OR_{event=disease} ≈ RR_{event=disease}

When the disease is rare, OR_{event=exposure} ≈ RR_{event=disease}