

## Chapter 3 - Probability

**Dice rolls.** (3.6, p. 92) If you roll a pair of fair dice, what is the probability of

(a) getting a sum of 1?

This is impossible, as all values on a fair die are integers. The probability of a sum of 1 is 0.

(b) getting a sum of 5?

There are  $6 \times 6 = 36$  possibilities for the sum of two independent rolls. There are 4 outcomes which result in a sum of 5. The probability is  $1/9$ .

(c) getting a sum of 12?

Only one combination will result in a sum of 12. The probability is  $1/36$ .

---

**Poverty and language.** (3.8, p. 93) The American Community Survey is an ongoing survey that provides data every year to give communities the current information they need to plan investments and services. The 2010 American Community Survey estimates that 14.6% of Americans live below the poverty line, 20.7% speak a language other than English (foreign language) at home, and 4.2% fall into both categories.

(a) Are living below the poverty line and speaking a foreign language at home disjoint?

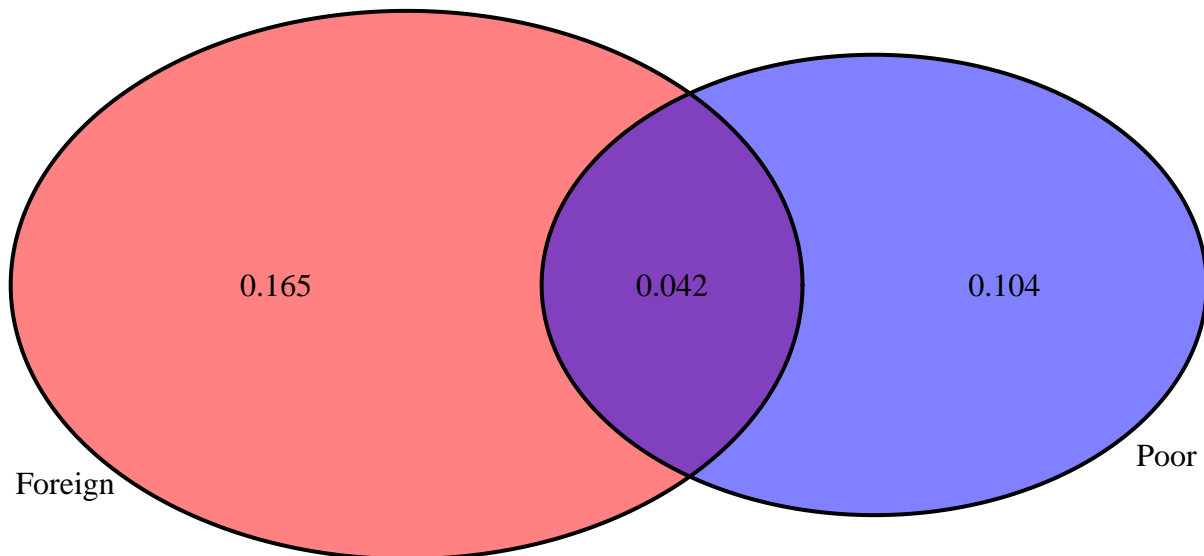
No. They are not mutually exclusive, a person can do both.

(b) Draw a Venn diagram summarizing the variables and their associated probabilities.

```
library(VennDiagram)

## Loading required package: grid
## Loading required package: futile.logger

grid.newpage()
draw.pairwise.venn(area1 = 0.146,
                   area2 = 0.207,
                   cross.area = 0.042,
                   fill = c("blue", "red"),
                   category = c("Poor", "Foreign"))
```



(c) What percent of Americans live below the poverty line and only speak English at home?

```
P_poor <- 0.146
P_foreign <- 0.207
P_foreignAndPoor <- 0.042

P_anglophoneAndPoor <- P_poor - P_foreignAndPoor
P_anglophoneAndPoor

## [1] 0.104
```

(d) What percent of Americans live below the poverty line or speak a foreign language at home?

```
P_foreignOrPoor <- P_poor + P_foreign - P_foreignAndPoor  
P_foreignOrPoor
```

```
## [1] 0.311
```

(e) What percent of Americans live above the poverty line and only speak English at home?

```
P_notPoorAndAnglophone <- (1 - P_poor) * (1 - P_foreign)  
P_notPoorAndAnglophone
```

```
## [1] 0.677222
```

(f) Is the event that someone lives below the poverty line independent of the event that the person speaks a foreign language at home?

Since the majority of Americans appear to live above the poverty line and speak only English at home, these do not appear to be independent variables.

---

**Assortative mating.** (3.18, p. 111) Assortative mating is a nonrandom mating pattern where individuals with similar genotypes and/or phenotypes mate with one another more frequently than what would be expected under a random mating pattern. Researchers studying this topic collected data on eye colors of 204 Scandinavian men and their female partners. The table below summarizes the results. For simplicity, we only include heterosexual relationships in this exercise.

		<i>Partner (female)</i>			Total
		Blue	Brown	Green	
<i>Self (male)</i>	Blue	78	23	13	114
	Brown	19	23	12	54
	Green	11	9	16	36
	Total	108	55	41	204

- (a) What is the probability that a randomly chosen male respondent or his partner has blue eyes?

$$(78 + 23 + 13 + 19 + 11) / 204 = 0.7058824$$

- (b) What is the probability that a randomly chosen male respondent with blue eyes has a partner with blue eyes?

$$78 / 114 = 0.6842105$$

- (c) What is the probability that a randomly chosen male respondent with brown eyes has a partner with blue eyes? What about the probability of a randomly chosen male respondent with green eyes having a partner with blue eyes?

$$19 / 54 = 0.3518519 \quad 11 / 36 = 0.3055556$$

- (d) Does it appear that the eye colors of male respondents and their partners are independent? Explain your reasoning.

The counts of respondents with same-colored eyes are clearly larger than any individual other combinations. Partners' eye colors do not appear to be independent.

**Books on a bookshelf.** (3.26, p. 114) The table below shows the distribution of books on a bookcase based on whether they are nonfiction or fiction and hardcover or paperback.

	<i>Format</i>		Total	
	Hardcover	Paperback		
<i>Type</i>	Fiction	13	59	72
	Nonfiction	15	8	23
	Total	28	67	95

- (a) Find the probability of drawing a hardcover book first then a paperback fiction book second when drawing without replacement.

$$(28/95) * (59/94) = 0.1849944$$

- (b) Determine the probability of drawing a fiction book first and then a hardcover book second, when drawing without replacement.

$$((59/95) * (28/94)) + ((13/95) * (27/94)) = 0.2243001$$

- (c) Calculate the probability of the scenario in part (b), except this time complete the calculations under the scenario where the first book is placed back on the bookcase before randomly drawing the second book.

$$(72/95) * (28/95) = 0.2233795$$

- (d) The final answers to parts (b) and (c) are very similar. Explain why this is the case.

Since there are nearly a hundred books on the shelf, replacing or not replacing for a single trial doesn't affect the probability too much. If there were 4 consecutive pulls without replacement, the answers would diverge more against 4 pulls with replacement.

**Baggage fees.** (3.34, p. 124) An airline charges the following baggage fees: \$25 for the first bag and \$35 for the second. Suppose 54% of passengers have no checked luggage, 34% have one piece of checked luggage and 12% have two pieces. We suppose a negligible portion of people check more than two bags.

- (a) Build a probability model, compute the average revenue per passenger, and compute the corresponding standard deviation.

```
probability <- c(0.54, 0.34, 0.12)
price <- c(0, 25, 60)

mu <- sum(price * probability)
dev <- sd(price * probability)

mu
```

```
## [1] 15.7
```

```
dev
```

```
## [1] 4.578573
```

- (b) About how much revenue should the airline expect for a flight of 120 passengers? With what standard deviation? Note any assumptions you make and if you think they are justified.

```
passengers <- 120

noquote(c("expected revenue:", mu * passengers))

## [1] expected revenue: 1884

noquote(c("standard deviation:", dev * passengers))

## [1] standard deviation: 549.428794294584

noquote("example simulation")

## [1] example simulation

set.seed(69420)
sum(sample(price, size = passengers,
           prob = probability, replace = TRUE))

## [1] 2070
```

---

**Income and gender.** (3.38, p. 128) The relative frequency table below displays the distribution of annual total personal income (in 2009 inflation-adjusted dollars) for a representative sample of 96,420,486 Americans. These data come from the American Community Survey for 2005-2009. This sample is comprised of 59% males and 41% females.

<i>Income</i>	<i>Total</i>
\$1 to \$9,999 or loss	2.2%
\$10,000 to \$14,999	4.7%
\$15,000 to \$24,999	15.8%
\$25,000 to \$34,999	18.3%
\$35,000 to \$49,999	21.2%
\$50,000 to \$64,999	13.9%
\$65,000 to \$74,999	5.8%
\$75,000 to \$99,999	8.4%
\$100,000 or more	9.7%

**(a) Describe the distribution of total personal income.**

The distribution is almost normal, with a fat right tail.

**(b) What is the probability that a randomly chosen US**

resident makes less than \$50,000 per year?

$$P(\text{less than 50k}) = 62.2\%$$

**(c) What is the probability that a randomly chosen US resident makes less than \$50,000 per year and is female?**

Note any assumptions you make.

Assuming these figures apply evenly to male and female participants:

$$P(\text{female and less than 50k}) = 0.25502 \cdot 0.622 \cdot 0.41 = 25.502\%$$

**(d) The same data source indicates that 71.8% of females make less than \$50,000 per year. Use this value to determine whether or not the assumption you made in part (c) is valid.**

$$P(\text{less than 50k} \mid \text{female}) = 0.718$$

$$\text{Expected } P(\text{less than 50k} \mid \text{female}) = P(\text{female and less than 50k}) / P(\text{female}) = 0.622$$

Of course our assumption is incorrect.....