

Lab 4: Screening and Precision Public Health

Your name and student ID

today's date

Screening

We will talk about a very important use of conditional probability in public health and medicine, which is the idea of tools that screen for health outcomes. There are many examples of this, including mammograms for detection of breast cancer, the prostate specific antigen (PSA) for detection of prostate cancer, as well as tests for exposure to infectious diseases, and so forth. We will consider two types of events: i) whether the subject truly has the health condition of interest (let D denote the disease of interest, and D' its complement), and ii) whether a test was positive or not for this outcome (let $+$ denote a positive test and $-$ denote its complement, a negative test). There are several statistics that are used to evaluate the performance of a test, some of which are derivable from each other:

- Sensitivity: $P(+ | D)$ or the probability of test being positive if one has the disease.
- Specificity: $P(- | D')$ or the probability of test being negative given one does not have the disease.
- Positive predictive value (PPV): $P(D | +)$ or the probability of having the disease if an individual tests positive.
- Negative predictive value (NPV): $P(D' | -)$ or the probability of not having the disease if an individual tests negative.

Consider the following situation: Assume the total number of subjects is 10,000 and that, $P(D) = 0.05$, $P(+ | D) = 0.95$, $P(- | D') = 0.95$. This set up implies that the disease is rare, but that a very accurate test exists (i.e., equally high sensitivity and specificity).

**1. Fill in the following two-way table with the absolute frequencies using the information provided in the problem. (You can fill it in simply by typing the frequencies in the nine empty cells.):

	D	D'	Total
$+$			
$-$			
Total			

Create a vector, `p1`, with the entries of this table, entering one column at a time. For example, if my table was as follows:

	D	D'	Total
$+$	1	2	3
$-$	4	5	6
Total	7	8	9

Then I would enter my vector as:

```
p1 <- c(1, 4, 7, 2, 5, 8, 3, 6, 9)
```

```
p1 <- "<<<<YOUR CODE HERE>>>>"
p1
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
check_problem1()
```

```
## [1] "Checkpoint 1 Error: You did not make a numeric vector."
## [1] "Checkpoint 2 Error: Your list has the wrong number of elements"
## [1] "Checkpoint 3 Error: Your first column has a wrong number."
## [1] "Checkpoint 4 Error: Your second column has a wrong number."
## [1] "Checkpoint 5 Error: Your row totals are wrong."
##
## Problem 1
## Checkpoints Passed: 0
## Checkpoints Errored: 5
## 0% passed
## -----
## Test: FAILED
```

**2. Calculate the PPV using the entries from your table. Store the PPV as a percentage.

```
p2 <- "<<<<YOUR ANSWER HERE>>>>"
```

```
check_problem2()
```

```
## [1] "Checkpoint 1 Error: Is p2 a numeric value?"
## [1] "Checkpoint 2 Error: Is your answer converted to percent?"
## [1] "Checkpoint 3 Error: Did you compute the correct probability?"
##
## Problem 2
## Checkpoints Passed: 0
## Checkpoints Errored: 3
## 0% passed
## -----
## Test: FAILED
```

****3.** Re-do the two-way table and re-calculate the PPV, this time assuming that $P(D) = 0.02$. (Note that $P(+ | D) = 0.95$, $P(- | D') = 0.95$, as with the previous question.)

	D	D'	Total
+			
-			
Total			

Similarly, create a vector, `p3`, with the entries of this table, entering one column at a time.

```
p3 <- "<<<<YOUR CODE HERE>>>>"
p3
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
check_problem3()
```

```
## [1] "Checkpoint 1 Error: You did not make a numeric vector."
## [1] "Checkpoint 2 Error: Your list has the wrong number of elements"
## [1] "Checkpoint 3 Error: Your first column has a wrong number."
## [1] "Checkpoint 4 Error: Your second column has a wrong number."
## [1] "Checkpoint 5 Error: Your row totals are wrong."
##
## Problem 3
## Checkpoints Passed: 0
## Checkpoints Errored: 5
## 0% passed
## -----
## Test: FAILED
```

**4. Again, calculate the PPV using the entries from your table. Store the PPV as a percentage rounded to 1 decimal place.

```
p4 <- "<<<<YOUR ANSWER HERE>>>>"
```

```
check_problem4()
```

```
## [1] "Checkpoint 1 Error: Is p2 a numeric value?"
## [1] "Checkpoint 2 Error: Is your answer converted to percent?"
## [1] "Checkpoint 3 Error: Did you round your answer?"
## [1] "Checkpoint 4 Error: Did you compute the correct probability?"
##
## Problem 4
## Checkpoints Passed: 0
## Checkpoints Errored: 4
## 0% passed
## -----
## Test: FAILED
```

****5.** Explain why the sensitivity is so high, but the PPV is low for the first calculation and even lower for the second calculation.

<TODO: YOUR ANSWER HERE>

Precision public health

One of the goals of public health research is to group people by risk factors or demographic variables so that decision-makers can predict, with actionable accuracy, which groups are at high and low risk of an adverse health outcome. In this set of questions, we consider stratified two-way tables, which are two-way tables specific to levels of a third grouping variable. Here, the adverse health outcome is coronary heart disease (CHD), which we represent by D . We study two categorical risk factors, smoking (defined by S for smoking and S' for no smoking) and age (defined by A for older age and A' for younger age).

First, read in the aggregated data set. The last column (n) is the number of individuals in each group.

```
library(dplyr)
library(tidyverse)
chd_dat <- read_csv("CHD.csv")

chd_dat
```

```
## # A tibble: 8 x 4
##   Age      Smoking CHD         n
##   <chr> <chr>   <chr> <dbl>
## 1 young yes     yes      60
## 2 young yes     no     240
## 3 young no      yes    105
## 4 young no      no    595
## 5 old  yes     yes    180
## 6 old  yes     no    120
## 7 old  no      yes    210
## 8 old  no      no    490
```

From this table:

6. Calculate the following probabilities. Convert your answers to percentages at the nearest whole number. Input the values (without the %) in the order below into the vector p6:

- $P(D \mid A', S)$
- $P(D \mid A', S')$
- $P(D \mid A, S)$
- $P(D \mid A, S')$

```
p6 <- c("YOUR ANSWER HERE")
p6
```

```
## [1] "YOUR ANSWER HERE"
```

```
check_problem6()
```

```
## [1] "Checkpoint 1 Error: You did not make a numeric vector."
## [1] "Checkpoint 2 Error: Your list needs 4 elements"
## [1] "Checkpoint 3 Error: Incorrect probability for P(D|A', S)."
## [1] "Checkpoint 4 Error: Incorrect probability for P(D|A', S')."
## [1] "Checkpoint 5 Error: Incorrect probability for P(D|A, S)"
## [1] "Checkpoint 6 Error: Incorrect probability for P(D|A, S)'"
```

```
##  
## Problem 6  
## Checkpoints Passed: 0  
## Checkpoints Errored: 6  
## 0% passed  
## -----  
## Test: FAILED
```

If you prefer, you can do these calculations by hand based on `chd_dat`. Some students might wish to use R commands to calculate these probabilities. There are **many** ways to do this. You could use `dplyr` functions to perform the calculations. Alternatively, here are some new functions for those of you interested in learning more R. (Note that these new functions won't be tested, they are for your information only, and to expose you to more of the R language!). You could consider using the `uncount()` function to expand the data based upon the numbers in each group (i.e., `n`) and assign the expanded data frame to a new name. Then, you can use the `tabyl` function from the `janitor` package to create stratified two-way tables, and the relevant `adorn_` functions from `janitor` to convert the frequencies to percentages.

****7.** What do these numbers imply about smoking and the risk of CHD?

<TODO: YOUR ANSWER HERE>

**8. Calculate the marginal probability of CHD. This can be written as $P(D)$. Store it in the variable p8.

```
p8 <- "YOUR ANSWER HERE"
p8
```

```
## [1] "YOUR ANSWER HERE"
```

```
check_problem8()
```

```
## [1] "Checkpoint 1 Error: Is p8 a numeric value?"
## [1] "Checkpoint 2 Passed: Correct!"
## [1] "Checkpoint 3 Error: Did you round your answer?"
## [1] "Checkpoint 4 Error: Did you compute the correct probability?"
##
## Problem 8
## Checkpoints Passed: 1
## Checkpoints Errored: 3
## 25% passed
## -----
## Test: FAILED
```

**9. Calculate the conditional probabilities The $P(D \mid A')$ and $P(D \mid A)$. Store them in the vector p9 rounded to the nearest whole number percentage.

```
p9 <- c("YOUR ANSWER HERE")
p9
```

```
## [1] "YOUR ANSWER HERE"
```

```
check_problem9()
```

```
## [1] "Checkpoint 1 Error: You did not make a numeric vector."
## [1] "Checkpoint 2 Error: Your list needs 2 elements"
## [1] "Checkpoint 3 Error: Incorrect probability for P(D|A')."
## [1] "Checkpoint 4 Error: Incorrect probability for P(D|A)"
##
## Problem 9
## Checkpoints Passed: 0
## Checkpoints Errored: 4
## 0% passed
## -----
## Test: FAILED
```

**10. If you had an intervention that could eliminate the risk of smoking on CHD, which group (defined by age) would see the biggest change in CHD from this intervention?

<TODO: YOUR ANSWER HERE>

Check your score

Click on the middle icon on the top right of this code chunk (with the downwards gray arrow and green bar) to run all your code in order. Then, run this chunk to check your score.

```
# Just run this chunk.  
total_score()
```

##		Test	Points_Possible	Type
## Problem 1		FAILED	1	autograded
## Problem 2		FAILED	1	autograded
## Problem 3		FAILED	1	autograded
## Problem 4		FAILED	1	autograded
## Problem 5	NOT YET GRADED		1	free-response
## Problem 6		FAILED	1	autograded
## Problem 7	NOT YET GRADED		1	free-response
## Problem 8		FAILED	1	autograded
## Problem 9		FAILED	1	autograded
## Problem 10	NOT YET GRADED		1	free-response

Submission

For assignments in this class, you'll be submitting using the **Terminal** tab in the pane below. In order for the submission to work properly, make sure that:

1. Any image files you add that are needed to knit the file are in the **src** folder and file paths are specified accordingly.
2. You **have not changed the file name** of the assignment.
3. The file is saved (the file name in the tab should be **black**, not red with an asterisk).
4. The file knits properly.

Once you have checked these items, you can proceed to submit your assignment.

1. Click on the **Terminal** tab in the pane below.
2. Copy-paste the following line of code into the terminal and press enter.

```
cd; cd ph142-sp20/lab/lab04; python3 turn_in.py
```

3. Follow the prompts to enter your Gradescope username and password. When entering your password, you won't see anything come up on the screen—don't worry! This is just for security purposes—just keep typing and hit enter.
4. If the submission is successful, you should see "Submission successful!" appear as output.
5. If the submission fails, try to diagnose the issue using the error messages—if you have problems, post on Piazza.

The late policy will be strictly enforced, **no matter the reason**, including submission issues, so be sure to submit early enough to have time to diagnose issues if problems arise.