

# Unit 4 – Categorical Data Analysis

## Practice Problems (2 of 2)

**Due date extended: Friday December 13, 2024**

### Before You Begin - R Users

\_\_\_1. Be sure you have downloaded the dataset **hersdata.Rdata**

\_\_\_2. The supplementary R illustration makes use of the following packages. Make sure you have done a one time installation of them, taking care to remember that R is case-sensitive.

**gmodels, DescTools, summarytools, tidyverse**

#1. *Source:* Rosenman RH, Friedman M, Straus R, Wurm M, Kositchek R, Hahn W and Werthessen NG (1964) A predictive study of coronary heart disease: the western collaborative group study. *Journal of the American Medical Association*, **189**, 113-120.

Note to class: This study was used in several data analysis illustrations in the book by Vittinghoff et al (Regression Methods in Biostatistics: Linear, Logistic, Survival and Repeated Measures Models, 2<sup>nd</sup> edition, Springer, 2012).

This exercise gives you practice performing a **2 x K test of trend** for contingency table data.

The Western Collaborative Group Study (WCGS) was a prospective study of 3,154 men, all initially disease-free, who were followed for events of coronary heart disease (CHD). At the eight years follow-up mark, there were 257 events of CHD. Several potential predictors of CHD were of interest, including: age (at enrollment), cholesterol, systolic blood pressure, hematocrit, ECG status, smoking, and relative weight.

Age is a continuous variable, suggesting that it be modelled as such. However, the best way to model age is a question. Should it be treated as a linear predictor? Or modeled using a polynomial? Or modeled in some other way? To address this issue, an appropriate preliminary analysis might group study participants into age groups (a discrete predictor that is ordinal!). Following are the data:

CHD	Age at Enrollment (years)					Total =
	35-40	41-45	46-50	51-55	55-60	
no	512	1036	680	463	206	2987
yes	31	55	70	65	36	257
<b>Total =</b>	543	1091	750	528	242	3154

1a. Perform a test of **general association**

H<sub>0</sub>: The odds of event of CHD is independent of interval of age

H<sub>A</sub>: The odds of event of CHD is associated with (differs with) interval of age

1b. Perform a test of **trend**

H<sub>0</sub>: The odds of event of CHD is independent of interval of age

H<sub>A</sub>: The odds of event of CHD is increases with interval of age

#2. Source: Triola MM and Triola MF. *Biostatistics for the Biological and Health Sciences* Boston: Pearson Addison Wesley, John Wiley, 2006. Chapter 10, Section 10-2. page 491.

This exercise gives you practice performing a [chi square goodness of fit \(GOF\) test](#).

Just briefly. Researchers suspect that people tend to self-report their weights lower than what they actually are and, in particular, tend to round down. If this is true, then one might expect to find that the last digits of self-reported weight are disproportionately often the digits “0”, “1”, “2”, “3”, “4” or “5”. The following table are the values of the last digit of self-reported weights for a sample of  $n=80$ . Carry out an appropriately chi square goodness-of-fit (GOF) test to assess if there is statistically significant evidence in this sample of the suspected phenomenon of “rounding down” when self-reporting weight.

	Last Digit of Self-Reported Weight									
	0	1	2	3	4	5	6	7	8	9
Frequency	35	0	2	1	4	24	1	4	7	2

2a. What are the null ( $H_0$ ) and alternative ( $H_A$ ) hypotheses?

2b. How many degrees of freedom does your chi square statistic have?

2c. What is the value of your chi square statistic?

2d. What is the p-value?

2e. In 1-2 sentences, what do you conclude?

### Supplement - Learn R (nothing to turn in)

Practice with missing values, creating a 0/1 variable and using the package {tidyverse}

## Introduction to The Heart and Estrogen/Progestin Replacement Study (HERS)

Source:

Hulley S, Grady D, Bush T, Furberg C, Herrington D, Riggs B and Vittinghoff E (1998). *Randomized trial of estrogen plus progestin for secondary prevention of heart disease in postmenopausal women. The Heart and Estrogen/progestin Replacement Study. Journal of the American Medical Association*, **280**(7), 605-613.

In the HERS study, Hulley et al. (1998) sought to determine if exercise, a modifiable behavior, might lower the risk of diabetes in non-diabetic women who were at risk of developing the disease. The question is a complex one because there are many risk factors for diabetes. Moreover, the type of woman who chooses to exercise may be related in other ways to risk of diabetes, apart from the fact of her exercise habit. For example, women who exercise regularly are typically younger and have lower body mass index (BMI); these characteristics also confer a risk benefit with respect to diabetes. Finally, the benefit of exercise may be mediated through a reduction of body mass index. Vittinghoff, Glidden, Shiboski and McCulloch (2005) consider portions of this data in their 2005 text, Regression Methods in Biostatistics: Linear, Logistic, Survival and Repeated Measures Models (Springer). Their dataset has  $n=2,763$  observations on 37 variables.

Here, we consider just 5 variables.

### Data Dictionary - Partial Listing

Position	Variable	Variable Label	Type	Codes	# missing (NA)
1	HT	Hormone Therapy	character	"hormone therapy" "placebo"	None
2	drinkany	Current drinker	character	"no" "yes"	2
3	LDL	LDL Cholesterol, mg/dl	numeric	Range: [36.8, 365.2]	11
4	SBP	Systolic, mm Hg	numeric	Range: [83.0, 224.0]	0
5	weight	Weight, kg	numeric	Range: [37.5, 132.0]	2

### Practice #1 -

**Reminder.** R will not analyze a character variable. If you want to treat it as a "categorical variable", you must create a factor version.

Load `hersdata.Rdata`. Use `class()` to check the datatype of the variable `drinkany`. Use `factor()` to convert `drinkany` to factor type. Check.

```
load(file="hersdata.Rdata")

class(hersdata$drinkany)
## [1] "character"

hersdata$drinkany <- factor(hersdata$drinkany)

class(hersdata$drinkany)
## [1] "factor"
```

### Practice #2 -

Use `library()` to attach the package `{summarytools}`. Use the function `freq()` in attached `{summarytools}` to produce a one way frequency table of `drinkany`. How many missing values are there?

```
library(summarytools)
freq(hersdata$drinkany)

## Frequencies
## hersdata$drinkany
## Type: Factor
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##      no    1680    60.848    60.848    60.803    60.803
##      yes    1081    39.152   100.000    39.124    99.928
##      <NA>      2      0.072    100.000    0.072   100.000
##      Total   2763   100.000   100.000   100.000   100.000
```

There are 2 missing values.

### Practice #3 -

**Important.** Whenever you create a new variable, always handle missing values explicitly (you'd be surprised what can happen when you don't!) Create a 0/1 numeric variable `drinkany01` that is coded as below; note that it handles missing values explicitly. Check:

**drinkany01** = 0 if **drinkany** == "no"  
 1 if **drinkany** == "yes" and  
 NA if **drinkany** == NA

```
hersdata$drinkany01 <- NA
hersdata$drinkany01[hersdata$drinkany=="no"] <- 0
hersdata$drinkany01[hersdata$drinkany=="yes"] <- 1

# Initialize to missing
# when drinkany=="no" assign 0 to drinkany01
# when drinkany=="yes" assign 1 to drinkany01

table(hersdata$drinkany, hersdata$drinkany01, useNA="always")
##
##          0    1 <NA>
## no    1680    0    0
## yes      0 1081    0
## <NA>      0    0    2
```

### Practice #4 -

The functions `filter()` and `select()` are in the package `{dplyr}` which is a core package in `{tidyverse}`. `{dplyr}` is attached automatically when you attach `{tidyverse}`. `filter()` and `select()` make subsetting data very easy!

Create a subset of `hersdata` called `mytiny` as follows:

Include only the following variables: **HT**, **LDL**, and **SBP**

Include only the observations with: **weight > 125**

Show

```
library(tidyverse)

mytiny <- hersdata %>%
  filter(weight > 125) %>%
  select(HT, LDL, SBP)

mytiny
```

*# use hersdata. THEN*  
*# filter( ) to choose observations to keep. AND THEN*  
*# select( ) to choose variables to keep*

		HT	LDL	SBP
## 1	hormone therapy	122.2	129	
## 2	placebo	204.6	133	
## 3	placebo	161.2	112	
## 4	placebo	137.0	130	
## 5	placebo	148.4	139	

Save.

```
save(mytiny, file="mytiny.Rdata")
```