

# STA 4504/5503 - Practice set 1 (with solutions)

## True or False

1. In  $2 \times 2$  tables, statistical independence is equivalent to a population odds ratio value of  $\theta = 1.0$ .
2. A British study reported in the New York Times: (Dec. 3, 1998) stated that of smokers who get lung cancer, “women were 1.7 times more vulnerable than men to get small-cell lung cancer.” The number 1.7 is a sample odds ratio.
3. Using data from the Harvard Physician’s Health Study, we find a 95% confidence interval for the relative risk relating having a heart attack to drug (placebo, aspirin) to be (1.4, 2.3). If we had formed the table with aspirin in the first row (instead of placebo), then the 95% confidence interval would have been  $(1/2.3, 1/1.4) = (.4, .7)$ .
4. Pearson’s chi-squared test of independence treats both the rows and the columns of the contingency table as nominal scale; thus, if either or both variables are ordinal, the test ignores that information.
5. For testing independence with random samples, Pearson’s  $X^2$  statistic and the likelihood-ratio  $G^2$  statistic both have exact chi-squared distributions for any sample size, as long as the sample was randomly selected.
6. Fisher’s exact test is a test of the null hypothesis of independence for  $2 \times 2$  contingency tables that fixes the row and column totals and uses a hypergeometric distribution for the count in the first cell. For a one-sided alternative of a positive association (i.e., odds ratio  $> 1$ ), the p-value is the sum of the probabilities of all those tables that have count in the first cell at least as large as observed, for the given marginal totals.
7. The difference of proportions, relative risk, and odds ratio are valid measures for summarizing  $2 \times 2$  tables for either prospective or retrospective (e.g., case-control) studies.
8. An ordinary regression model that treats the response  $Y$  as having a normal distribution is a case of a generalized linear model, with normal (a.k.a. Gaussian) random component, identity link function and assuming the same predictors, i.e. same systematic component.

## Open Ended Problems

- Each of 100 multiple-choice questions on an exam has five possible answers but only one correct response. For each question, a student randomly selects one response as the answer.
  - Specify the probability distribution of the student's number of correct answers on the exam, identifying the parameter(s) for that distribution.
  - Would it be surprising if the student made at least 50 correct responses? Explain your reasoning.
- Consider the following data from a women's health study (MI is myocardial infarction, i.e., heart attack).

|                     |            | MI  |     |
|---------------------|------------|-----|-----|
|                     |            | Yes | No  |
| Oral Contraceptives | Used       | 23  | 34  |
|                     | Never Used | 35  | 132 |

- Construct a 95% confidence interval for the population odds ratio.
- Does it seem plausible that the variables are independent? Explain.

3. For adults who sailed on the Titanic on its fateful voyage, the odds ratio between gender (female, male) and survival (yes, no) was 11.4.
  - (a) What is wrong with the interpretation, “The probability of survival for females was 11.4 times that for males.”
  - (b) When would the quoted interpretation be approximately correct? Why?
  - (c) The odds of survival for females equaled 2.9. For each gender, find the proportion who survived.
4. Explain two ways in which the generalized linear model extends the ordinary regression model that is commonly used for quantitative response variables.

5. In a recent General Social Survey, gender was cross-classified with party identification. The output below shows some results.

```
> gp
      party
gender dem indep rep
female 279    73 225
male   165    47 191

> addmargins(gp)
      party
gender dem indep rep Sum
female 279    73 225 577
male   165    47 191 403
Sum     444   120 416 980

> gp.chisq <- chisq.test(gp)
> gp.chisq
Pearson's Chi-squared test
data: gp
X-squared = 7.0095, df = 2, p-value = 0.03005

> gp.chisq$expected
      party
gender dem    indep    rep
female 261.42 70.653 244.93
male   182.58 49.347 171.07

> round(myadjresids(gp), 2)
      party
gender dem  indep  rep
female 2.29  0.46 -2.62
male  -2.29 -0.46  2.62
```

- (a) Explain what the numbers in the “expected” table represent. Show how to obtain 261.42.

- (b) Explain how to interpret the p-value given for the Chi-square statistic.
  - (c) Explain how to interpret the output of the last command (standardized adjusted residuals). Which counts were significantly higher than one would expect if party identification were independent of gender?
6. From a simulated data set containing information on 6047 customers such as whether the customer defaulted, is a student, and the average balance carried by the customer (ranging from 0 to 2700).

```
> model1 <- glm(default ~ balance, family = "binomial", data = train)
> summary(model1)
```

Coefficients:

|             | Estimate   | Std. Error | z value | Pr(> z )   |
|-------------|------------|------------|---------|------------|
| (Intercept) | -1.101e+01 | 4.887e-01  | -22.52  | <2e-16 *** |
| balance     | 5.669e-03  | 2.949e-04  | 19.22   | <2e-16 *** |

---

Null deviance: 1723.03 on 6046 degrees of freedom  
 Residual deviance: 908.69 on 6045 degrees of freedom  
 AIC: 912.69

- (a) Write out the model for predicting the probability of default and estimate the probability for a customer with a balance of 973.
- (b) Test whether balance is a significance predictor. Note there are two ways.
- (c) Using this output is there evidence of overdispersion? Explain.
- (d) Can we perform a goodness of fit test? True or False?