# STAT 518 --- Chapter 4 --- Contingency Tables

• Contingency tables are summaries (in matrix form) of categorical data, where the entries in the table are counts of how many observations fell into specific categories (and combinations of categories).

• A <u>one-way</u> contingency table summarizes data on a single categorical variable and has only one row.

• A <u>two-way</u> contingency table summarizes data on two categorical variables and may have several rows and several columns.

• Data on several categorical variables can be summarized by <u>multi-way</u> contingency tables.

• We begin with another goodness-of-fit test.

## Section 4.5:  Chi-Squared Goodness-of-Fit Test

• Suppose we have a single categorical variable with $c$ categories.  The cell counts can be arranged in a <u>one-way table</u>.

Example 1:  95 adults were randomly sampled and surveyed about their favorite sport.  There were 6 categories. Their preferences are summarized:

<u>Favorite Sport</u>

| Football | Baseball | Basketball | Auto | Golf | Other | $N$ |
|---|---|---|---|---|---|---|
| 37 | 12 | 17 | 8 | 5 | 16 | 95 |

$p_1$ = proportion of U.S. adults favoring football
$p_2$ = proportion of U.S. adults favoring baseball
$p_3$ = proportion of U.S. adults favoring basketball
$p_4$ = proportion of U.S. adults favoring auto racing
$p_5$ = proportion of U.S. adults favoring golf
$p_6$ = proportion of U.S. adults favoring "other"

• It was hypothesized that the true proportions are
$(p_1, p_2, p_3, p_4, p_5, p_6) = (.4, .1, .2, .06, .06, .18)$.

• We test our null hypothesis with the chi-squared goodness-of-fit test:

$H_0$:
$H_1$:  at least one of the hypothesized probabilities is wrong

The test statistic is:

where $O_j$ is the observed "cell count" for category $j$ and $E_j$ is the expected cell count for category $j$ if _____.

• Under $H_0$, $T$ has an asymptotic $\chi^2$ distribution with $c - 1$ d.f.

Decision Rule:

(large values of $T \rightarrow$ observed counts are very different from the expected counts under $H_0$)

<u>Assumptions</u>: (1) The data are at least <u>nominal</u>. (2) The random sample is sufficiently large. Koehler and Larntz's Rule of Thumb:

• If $H_0$ is true, expected cell count $E_j =$

**Example 1 data:**

| j | | | | |
|---|---|---|---|---|
| $O_j$ | | | | |
| $E_j$ | | | | |

**Test statistic value:**

**Decision Rule:**

**P-value**

**Conclusion:**

• See `chisq.test` function in R to perform this test.

# Chi-Squared Test with Unknown Parameters

• **If our null hypothesis specifies the distribution <u>except</u> for a certain number (say, *k*) of unknown parameters, we can adjust the chi-squared test to account for this.**

• **The main difference is that when *k* unknown parameters are estimated from the data, the asymptotic null distribution of *T* is $\chi^2$ with                 d.f.**

• **The unknown parameters must be estimated using "good methods" (see pp. 243-245):  Typically the method of moments or maximum likelihood estimators work well.**

**Example 2:  Page 244 lists data for the number of hits of 18 baseball players in their first 45 times at bat.  Is it reasonable that these data all follow the same binomial distribution with *n* = 45 and some unspecified *p*?**

• **To estimate the unknown *p*, we use the estimate:**

• **The expected cell counts can be found by the formula:**

• Note that some $E_j$ are very small; to alleviate this we should combine cells:

Test statistic value:

Decision Rule:

P-value

Conclusion:

• While contingency tables describe <u>discrete</u> data, the chi-squared test can be used to check goodness of fit for <u>continuous</u> models as well.

• In that case, the continuous data must be discretized by grouping into intervals.

• How to form the intervals is somewhat arbitrary.

**Example 1 from Section 6.2:  The data on page 445 consist of 50 observations.  At $\alpha = 0.05$, is it reasonable to claim that the data follow a normal distribution?**

**We first estimate the two unknown parameters ($\mu$ and $\sigma$) of the normal distribution:**

**Let's choose 5 intervals:**

| Interval | | | | | |
|---|---|---|---|---|---|
| $O_j$ | | | | | |
| $E_j$ | | | | | |

**Test statistic value:**

**Decision Rule:**

**P-value**

**Conclusion:**

# Section 4.1:  Tests for $2 \times 2$ Tables

• **Consider the simplest form of two-way table:**

• **Such a table could summarize data arising from**

  - **Having a single sample in which _____ _____ are measured on each individual**
  - **Having two samples in which _____ _____ is measured on each individual in each sample.**

### Comparing Two Probabilities, Independent Samples

• **Suppose we have two independent samples, with respective sizes $n_1$ and $n_2$.  We classify each individual in each sample into class 1 or class 2.**

• **Our data could be arranged in a $2 \times 2$ table as follows:**

• **The total number of observations is $N = n_1 + n_2$.**

• **Our goal is to compare the probability of "success" (Class 1) across the two populations:**

**Hypotheses:**

## Development of the Test Statistic

**As estimators of $p_1$ and $p_2$, we have:**

• **This estimates how far apart $p_1$ and $p_2$ are.**

• **Scaling this by dividing by the estimated standard error (see Eq. 5, p. 187), we get the test statistic**

**which has a _____ distribution when $H_0$ is true.**

• If $T_1$ is far from zero, this indicates that

• If $T_1$ is far below zero, this indicates that

• If $T_1$ is far above zero, this indicates that

**Decision Rules**

**H₁:**                               **H₁:**                               **H₁:**

**P-value:**

• **Note:  The normal approximation for $T_1$ is valid for large samples, say, if**

**Example 1:  A survey was conducted of 160 rural households and 261 urban households with Christmas trees.  Of interest was whether the tree was natural or artificial.  Is the probability of natural trees different for rural and urban households?  Use $\alpha = 0.05$.**

**Data:**

| Population | Tree | |
|---|---|---|
| | Natural | Artificial |
| Rural | 64 | 96 |
| Urban | 89 | 172 |

**H₀:**                      **H₁:**

**Test statistic:**

**Example 2:  Page 184 gives data from a study to determine whether a new lighting system worsened midshipmen's vision.**

**Data:**

| | | Vision | |
|---|---|---|---|
| | | Good | Poor |
| | Old | 714 | 111 |
| Lighting | | | |
| | New | 662 | 154 |

**H₀:**                      **H₁:**

**Test statistic:**

# Fisher's Exact Test

• In the previous examples, the row totals were the sizes of the two samples, which are <u>fixed</u> before the data are examined (i.e., they are not random).

• When we have a single sample in which two variables are measured on each individual, the resulting $2 \times 2$ table has _____ row totals and _____ column totals.
• We will cover that scenario in Section 4.2.

• In other situations, both the row totals and the column totals may be _____ prior to the data being examined.

• In this case of "_____ margins", Fisher's Exact Test is ideal.

Data setup:

• We again wish to compare:

Test statistic $T_2 =$

# Null Distribution

• Let $p$ = probability an observation is in Column 1.
• Under $H_0$, this probability is the same whether the observation is in Row 1 or Row 2.  Then:

P(table results | row totals) =

P(column totals) =

→ P(table results | row totals & column totals) =

• The decision is based on the P-value, which is found differently depending on the alternative hypothesis:

$H_1$:                          $H_1$:                          $H_1$:

• In all cases, reject $H_0$ if the p-value $\leq \alpha$.

Example 3:  Fourteen new hires (10 male and 4 female) are being assigned to bank positions (there are 4 account representative positions open and 10 (less desirable) teller positions open.  The data on page 190 summarize the assignments.  If all new employees are equally qualified, is there evidence that female hires were more likely to get the account representative jobs?

**Data:**


**H$_0$:**                              **H$_1$:**


**Test statistic:**
**P-value:**




• See `fisher.test` function in R to perform this test.

• Fisher's Exact Test may be used if the row totals and/or column totals are random, but in this case it is _____ _____ than the z-test.

• Fisher's Exact Test can also be viewed as an alternative to the z-test when the large-sample rule is not met, but the Exact Test _____ _____ when the sample size is very small.

• Suppose we have several related (but not identical) conditions in which sub-experiments are conducted, each of which produces a 2 × 2 table.

• It is of interest to see whether rows and columns are independent in <u>each</u> table.

# Mantel-Haenszel Test

• We assume we have $k \geq 2$ such $2 \times 2$ tables, each with fixed row and column totals (although the test can be done even with random totals).

Let $p_{1i} =$

and $p_{2i} =$

**Hypotheses:**

**Test statistic**

• The null distribution is approximately standard normal, tabulated in Table A1.

**Decision Rules and P-value:**




**Example 4: Three groups of cancer patients were given either a drug treatment or a control, and for each patient, whether the outcome was successful was recorded. Is there evidence that <u>in at least one group,</u> the treatment produces a better chance of success than the control? (Use $\alpha = 0.05$.)**

**Data:**




**H$_0$:**                                    **H$_1$:**




**Test statistic:**
**P-value:**

**Conclusion:**


• See `mantelhaen.test` function in R to perform this test.