

Week Five

Last Week

- Contingency Tables
 - Contingency Table Overview: multiple categorical variables
 - joint, marginal, and conditional probabilities
 - Comparing Proportions: Relative Risk & Odds Ratios
 - Chi-Squared Tests for Independence

This Week: Contingency Tables, Part II

Today:

- Activity
- Thursday: Lab

Next Week: Generalized Linear Models

Ordinal Data

Up to this point, we have largely considered nominal categorical data. However, we will also see ordinal data.

Table 1: Ordinal contingency table

| | low | medium | high |
|---------|-----------|-----------|-----------|
| group 1 | $n_{1,l}$ | $n_{1,m}$ | $n_{1,h}$ |
| group 2 | $n_{2,l}$ | $n_{2,m}$ | $n_{2,h}$ |

All procedures still apply in this setting. However, we might be interested in a different question rather than just independence we might wish to think more about the trends across the ordinal categories.

Consider simulating data with this structure.

1. Simulate ordinal responses with three levels across two different groups.
2. Visualize the data
3. Create a contingency table
4. Run a χ^2 test and interpret the results. Clearly articulate what you are testing here (avoid statistical lingo).
5. Create another table that includes the observed values and expected values (from the χ^2 test). How do these values contribute to the test statistic in the previous part?

The previous analysis treats the data as nominal, and is still valid, but doesn't explicitly account for the ordinal structure of the data. We can do something that mimics correlation with continuous data.

- First, we need to define scores for the each of the responses. An example would be low = 1, medium = 2, high = 3
- Then define $\bar{u} = \sum_i u_i p_{i+}$ to denote the sample mean of the row scores and $\bar{v} = \sum_j v_j p_{+j}$ to be the sample mean of the column scores.

$$r = \frac{\sum_{ij} (u_i - \bar{u})(v_j - \bar{v})p_{ij}}{\sqrt{[\sum_i (u_i - \bar{u})^2 p_{i+}] [\sum_j (v_j - \bar{v})^2 p_{+j}]}}$$

- $M^2 = (n-1)r^2$ and M^2 has an approximate χ^2 distribution with 1 degree of freedom (with large n).

Consider the following table.

| group | L | M | H |
|-------|-----|-----|-----|
| 1 | 298 | 161 | 126 |
| 2 | 131 | 163 | 121 |

calculate r and M^2 .

- r is

- $M^2 =$

Fisher's Exact Test

The story behind Fisher's exact test is that a lab mate claimed they could taste the difference between cups of tea that had milk added first or last. So an experiment was designed with 4 cups of each. We will do something similar on Thursday.

This setting requires the use of the hypergeometric distribution. Let n_{ij} be the cell counts for the i^{th} row and j^{th} column.

| | Milk(guess) | Tea (guess) | Total |
|--------------|-------------|-------------|----------|
| Milk(actual) | n_{11} | n_{12} | n_{1+} |
| Tea (guess) | n_{21} | n_{22} | n_{2+} |
| Total | n_{+1} | n_{+2} | |

Given we know the marginal counts, the guesses, n_{11} completely determines the other values (a bit like sudoku).

Using a hypergeometric distribution we can estimate the probability of the possible values for n_{11} as

$$Pr(n_{11}) = \frac{\binom{n_{1+}}{n_{11}} \binom{n_{2+}}{n_{+1}-n_{11}}}{\binom{n}{n_{+1}}}$$

```
tibble(n11 = c(0, 1, 2, 3, 4),
  Prob = c(choose(4,0) * choose(4,4) / choose(8,4),
    choose(4,1) * choose(4,3) / choose(8,4),
    choose(4,2) * choose(4,2) / choose(8,4),
    choose(4,3) * choose(4,1) / choose(8,4),
    choose(4,4) * choose(4,0) / choose(8,4))
) |> kable()
```

| n11 | Prob |
|-----|-----------|
| 0 | 0.0142857 |
| 1 | 0.2285714 |
| 2 | 0.5142857 |
| 3 | 0.2285714 |
| 4 | 0.0142857 |

How do we get p-values from this table?