# BIOS 6110
# Applied Categorical Data Analysis

Instructor: Yuan Huang, Ph.D.

Department of Biostatistics

Fall 2017

# Part VIII

# Models for Matched Pairs

## Overview

- Previously, we have worked on the data where the subjects are independently collected.

- This section introduces methods for comparing categorical responses for two samples that have a natural pairing between each subject in one sample and a subject in the other sample.

- Because each observation in one sample pairs with an observation in the other sample, the responses in the two samples are matched pairs.

- Because of the matching, the samples are statistically dependent.

- In this case, methods that treat the two sets of observations as independent samples are inappropriate

We will focus on 1-1 Match with binary response in this chapter.

## Objectives

1. Data display
   - Population-Average table
   - Subject-Specific tables

2. Compare two dependent proportions
   - McNemar test for testing marginal homogeneity
   - Estimate and CI of the proportion difference
   - Mantel-Haenszel odds ratio

3. Logistic regression for matched pairs
   - Marginal models
   - Conditional Logistic regression for matched-pairs
   - *Why unconditional logistic model is biased for matched-pairs

Two motivating data examples:

- 2000 General Social Survey (used in objective 2)

- Diabetes-MI Matched Pairs Case-Control Study (used in objective 3)

## Motivating Example: 2000 General Social Survey

In the 2000 General Social Survey, 1144 subjects were asked whether, to help the environment, they would be willing to (1) pay higher taxes or (2) accept a cut in living standards.

|                  | Cut Living Standards | | |
| ---------------- | --- | --- | ----- |
| Pay Higher Taxes | Yes | No  | Total |
| Yes              | 227 | 132 | 359   |
| No               | 107 | 678 | 785   |
| Total            | 334 | 810 | 1144  |

How can we compare the probabilities of a "yes" outcome for the two environmental questions?

## Comparing Dependent Proportions

Generally, we have

|                   | Cut Living Standards |         |          |
|-------------------|----------------------|---------|----------|
| Pay Higher Taxes  | Yes                  | No      | Total    |
| Yes               | $n_{11}$             | $n_{12}$ | $n_{1+}$ |
| No                | $n_{21}$             | $n_{22}$ | $n_{2+}$ |
| Total             | $n_{+1}$             | $n_{+2}$ | $n$      |

Define sample estimates

$$p_{ij} = n_{ij}/n, \quad i = 1, 2, j = 1, 2$$
$$p_{i+} = p_{i1} + p_{i2}, \quad i = 1, 2$$
$$p_{+j} = p_{1j} + p_{2j}, \quad j = 1, 2$$

The population value corresponding to $p_{ij}$ is denoted by $\pi_{ij}$. When $\pi_{+1} = \pi_{1+}$ and $\pi_{+2} = \pi_{2+}$, there is *marginal homogeneity*. From this condition, we have

$$\pi_{1+} - \pi_{+1} = (\pi_{11} + \pi_{12}) - (\pi_{11} + \pi_{21}) = \pi_{12} - \pi_{21}$$

That is, marginal homogeneity is equivalent to $\pi_{12} = \pi_{21}$.

## Testing Marginal Homogeneity: McNemar Test

For matched pairs data with a binary response, a test of marginal homogeneity has null hypothesis

$$H_0 : \pi_{12} = \pi_{21}, \quad \text{or equivalently} \quad H_0 : \pi_{+1} = \pi_{1+}$$

As off-diagonal counts are $n_{12} + n_{21}$, under $H_0$, both $n_{12}$ and $n_{21}$ are distributed as Binomial($n_{12} + n_{21}$, 0.5).

Construct test based on $n_{12}$. When $n_{12} + n_{21}$ is large, the binomial can be approximately by normal distribution with mean $(n_{12} + n_{21}) \times 0.5$ and variance $(n_{12} + n_{21}) \times (0.5) \times (1 - 0.5)$. Therefore, the z statistics can be constructed accordingly.

$$z = \frac{n_{12} - (n_{12} + n_{21}) \times 0.5}{\sqrt{(n_{12} + n_{21}) \times 0.5 \times 0.5}} = \frac{n_{12} - n_{21}}{\sqrt{n_{12} + n_{21}}}$$

The $z \sim \text{N}(0,1)$, or $z^2$ is approximately $\chi^2$ with df = 1 under $H_0$.

The chi-squared test ($z^2$) for a comparison of two dependent proportions is called McNemar test.

Example: 2000 General Social Survey (cont.)

The McNemar's test statistic is

$$\frac{(132 - 107)^2}{132 + 107} = 2.615$$

with *p*-value

```
> 1-pchisq(2.615,1)
[1] 0.1058575
```

There is no sufficient evidence that the probability of approval was greater for higher taxes than for a lower standard of living

## Estimating Differences of Proportions

The point estimate of $\pi_{1+} - \pi_{+1}$ is $p_{1+} - p_{+1}$

The estimated standard error (SE) of $p_{1+} - p_{+1}$ is

$$
\left\{ \frac{p_{1+}(1 - p_{1+}) + p_{+1}(1 - p_{+1}) - 2(p_{11}p_{22} - p_{12}p_{21})}{n} \right\}^{1/2}
$$
$$
= n^{-1}\sqrt{(n_{12} + n_{21}) - n^{-1}(n_{12} - n_{21})^2}
$$

95% Confidence interval: $(p_{1+} - p_{+1}) \pm 1.96 \times SE$

- Matched-pairs data usually show a positive association, which implies a odds ratio $> 1$. A sample odds ratio exceeding 1.0 corresponds to $p_{11}p_{22} > p_{12}p_{21}$, a negative contribution from the third term.

- Thus, an advantage of using dependent samples, rather than independent samples, is a smaller variance for the estimated difference in proportions.

Example: 2000 General Social Survey (cont.)

The point estimate is

$$359/1144 - 334/1144 = 0.314 - 0.292 = 0.022$$

The estimated SE is

$$1144^{-1}\sqrt{(132 + 107) - 1144^{-1}(132 - 107)^2} = 0.0135$$

A 95% CI is

$$0.022 \pm 1.96(0.0135) = (-0.004, 0.048)$$

## Population-Average Table vs Subject-Specific Tables

In the 2000 Social Survey Example, the previous table cross-classifies in a single table the two response for **all subjects**. This is called the **population-average table**.

We can consider **each matched pair as a cluster**, and present $2 \times 2$ partial table for each pair. In this case, each of the $2 \times 2$ partial table shows 2 observations from the matched pairs. Because there are 1144 subject, we will get $2 \times 2 \times 1144$ three-way table. This is called the **subject-specific tables**.

E.g., each partial table can be displayed as

|                       | Response |      |
| --------------------- | -------- | ---- |
|                       | Yes      | No   |
| Pay Higher Taxes      | XXX      | XXX  |
| Cut Living Standards  | XXX      | XXX  |

## Example: 2000 General Social Survey (cont.)

| (a) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 1 | 0 |
| Cut Living Standards | 1 | 0 |

| (b) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 1 | 0 |
| Cut Living Standards | 0 | 1 |

| (c) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 0 | 1 |
| Cut Living Standards | 1 | 0 |

| (d) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 0 | 1 |
| Cut Living Standards | 0 | 1 |

- There are 227 such partial table (a), which correspond to $n_{11}$.

- There are 132 such partial table (b), which correspond to $n_{12}$.

- There are 107 such partial table (c), which correspond to $n_{21}$.

- There are 678 such partial table (d), which correspond to $n_{22}$.

## Mantel-Haenszel Odds Ratio - Based on Subject-Specific Tables

Previously in XYZ three-way table, where $X$ is the binary treatment, $Y$ is binary response, and $Z$ is the centers ($k = 1, \ldots, K$).

|           | Diseased  | Non-diseased | Totals    |
|-----------|-----------|--------------|-----------|
| Exposed   | $n_{11k}$ | $n_{12k}$    | $n_{1+k}$ |
| Unexposed | $n_{21k}$ | $n_{22k}$    | $n_{2+k}$ |
| Totals    | $n_{+1k}$ | $n_{+2k}$    | $n_{++k}$ |

The Mantel-Haenszel method

- assumes that there is a true odds ratio which is consistent across $k$
- provides a pooled estimate of the common odds ratio. In essence, it is a weighted average of the odds ratios from the individual strata

The Mantel-Haenszel estimate of the odds ratio is

$$\hat{\theta}_{MH} = \frac{\sum_{k=1}^{K} n_{11k} n_{22k} / n_{++k}}{\sum_{k=1}^{K} n_{21k} n_{12k} / n_{++k}}$$

Note that, in our current matched pair case, the $K$ is number of pairs.

## Example: 2000 General Social Survey (cont.)

| (a) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 1 | 0 |
| Cut Living Standards | 1 | 0 |

| (b) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 1 | 0 |
| Cut Living Standards | 0 | 1 |

| (c) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 0 | 1 |
| Cut Living Standards | 1 | 0 |

| (d) | Response | |
|---|---|---|
| | Yes | No |
| Pay Higher Taxes | 0 | 1 |
| Cut Living Standards | 0 | 1 |

- Pairs as (a) $\rightarrow$ contribute 0 to denominator and 0 to numerator

- Pairs as (b) $\rightarrow$ contribute 0 to denominator and 1 to numerator

- Pairs as (c) $\rightarrow$ contribute 1 to denominator and 0 to numerator

- Pairs as (d) $\rightarrow$ contribute 0 to denominator and 0 to numerator

- $$\hat{\theta}_{MH} = \frac{\sum_{k=1}^{K} n_{11k} n_{22k} / n_{++k}}{\sum_{k=1}^{K} n_{21k} n_{12k} / n_{++k}} = \frac{n_{12}}{n_{21}} = \frac{132}{107} = 1.234$$

## Marginal Models for Marginal Proportions

A marginal table can be obtained by adding partial tables across stratas.

| Issue | Response | | |
|---|---|---|---|
| | Yes | No | Total |
| Pay Higher Taxes | 359 | 785 | 1144 |
| Cut Living Standards | 334 | 810 | 1144 |

A marginal model aims to study the marginal distributions of response for the two observations. The third dimension, i.e., pair stratas, is omitted.

$$\frac{359/785}{334/810} = 1.11.$$

The population odds of willingness to pay higher taxes are estimated to be 11% higher than the population odds of willingness to accept cuts in living standards.

## Marginal Models for Marginal Proportions (cont.)

- We already know that this analysis can also be obtained by using the following logistic regression

$$\text{logit}[\Pr(Y = \text{``Yes''})] = \alpha + \beta x$$

where $x$ is an indicator for question 1.

- Chapter 9 will discuss marginal models in details.