Compare Binary Data Between two groups

Review

Compare means

Compare medians

Setting	Parametric Test	Nonparametric Test
Matched pairs	Paired t-test (one sample t-test on the differences)	Wilcoxon signed rank test
Two independent samples	Two-sample t-test	Wilcoxon Rank Sum test

Data File

ID	group	outcome
1	1	611
2	1	621
3	1	614
4	1	593
5 6	1	593
6	1	653
7	1	600
8	1	554
9	1	603
10	1	569
11	2	635
12	2	605
13	2	638
14	2	594
15	2	599
16	2 2	632
17	2	631
18	2	588
19	2	607

Review

Compare means Compare medians **Nønparametric Parametric Setting Test lest** Paired itest Wilcoxon signed rank Matched pairs (one sample t-test test on the differences) Two Two-sample t-test Wilcoxon Rank Sum test independent samples

Chi-square Test or Fisher exact Test

Data File

ID	group	outcome
1	1	1
2	1	1
3	1	0
1 2 3 4	1	0
5 6	1	0
6	1	1
7	1	0
8	1	1
9	1	0
10	1	1
11	2	0
12	2	0
13	2	0
14	2	1
15	2	0
16	2	0
17	2 2 2 2 2 2 2	1
18	2	1
19	2	0

Risk

Goal: Compare risk of PVD in smokers and non-smokers?

Example: cross-sectional study of peripheral vascular disease (PVD) in Scottish men

	PVD		
Cigarette Smoker	Yes	No	Total
Yes	15	1712	1727
No	41	3188	3229
Total	56	4900	4956

Risk is the Prob(disease/exposure to risk factor)

Measure Association of Exposure and Disease

Risk Difference

Risk Ratio (Relative Risk)



Odds Ratio



Calculate Risks

	PVD			
Cigarette Smoker	Yes	No	Total	
Yes	15	1712	1727	
No	41	3188	3229	
Total	56	4900	4956	

$$p_{\scriptscriptstyle S} \quad = \quad \text{probability of PVD in smokers}$$

= risk of PVD in smokers

$$\hat{p}_S = 15/1727 = 0.009$$

$$p_{\scriptscriptstyle NS}$$
 = probability of PVD in non-smokers

= risk of PVD in non-smokers

$$\hat{p}_{NS} = 41/3229 = 0.013$$

Risk Difference

Risk Difference

The magnitude of the difference (exposed vs. non-exposed) in risk often is interpreted in relation to the "baseline" (non-exposed) risk

Risk Ratio

• Risk ratio (relative risk)

$$\lambda = \frac{\text{risk of disease in exposed subjects}}{\text{risk of disease in non-exposed subjects}}$$

	PVD			
Cigarette Smoker	Yes	No	Total	
Yes	15	1712	1727	
No	41	3188	3229	
Total	56	4900	4956	

$$\lambda = \frac{p_S}{p_{NS}} = \frac{15/1727}{41/3229} = 0.68$$

Interpretation: We estimate that the risk (probability) of PVD in smokers is 68% of the risk of PVD in non-smokers

Odds

- Risk is a probability that quantifies the likelihood that an event occurs
- We may instead be interested in the likelihood of an event in relation to how unlikely it is
- Such a quantity is known as the odds of an event:

$$\mathsf{Odds} = \frac{\mathsf{Risk}}{1 - \mathsf{Risk}}$$

• Example: If a disease occurs in 25% of adults, we can also say adults have an odds of 1/3 of developing the disease.

Odds Ratio

- As relative risk is used to compare the risk of two groups, we can use the relative odds to compare the odds of two groups
- Relative odds is more commonly referred to as the odds ratio:

$$\psi = \frac{\mathsf{Odds} \; \mathsf{for} \; \mathsf{Group} \; 1}{\mathsf{Odds} \; \mathsf{for} \; \mathsf{Group} \; 2}$$

- An odds ratio is much harder to interpret than a relative risk
- However, an odds ratio is used much more often than a relative risk

Computing Odds Ratio from a 2 by 2 Table

• Thus, for a general 2 x 2 table of exposure and disease

	Dise		
Risk Factor	Yes	No	Total
Exposed	а	b	a+b
Non-exposed	С	d	c+d
Total	a+c	b+d	n

the odds ratio of disease for exposed versus non-exposed is

$$\widehat{\psi} = \frac{ad}{bc} \qquad \widehat{\psi} = \frac{\frac{a}{a+b} / \frac{b}{a+b}}{\frac{c}{c+d} / \frac{d}{c+d}}$$

$$= \frac{a/b}{d/d}$$

$$= \frac{ad}{bc}$$

Example

	PVD		
Cigarette Smoker	Yes	No	Total
Yes	15	1712	1727
No	41	3188	3229
Total	56	4900	4956

• Risk difference = 15/1727 - 41/3229 = -0.004

• Relative risk =
$$\frac{15/1727}{41/3229}$$
 = 0.6840

• Odds Ratio =
$$\frac{15 \times 3188}{41 \times 1712}$$
 = 0.68127

Chi-squared & Fisher's Exact Test

Test Association of Exposure and Disease

2 × 2 table			
F1	Disease		
Exposed	Yes	No	
Yes			
No			

Chi-squared Test

We can also view our hypothesis test as:

 H_o : $\psi = 1 \rightarrow$ disease is not associated with exposure

 H_a : $\psi \neq 1$ \to disease is associated with exposure Relative Risk or Odd Ratio

or

 H_o : $\psi = 1 \rightarrow$ disease and exposure are independent

 H_a : $\psi \neq 1$ \rightarrow disease and exposure are dependent

 We can compare these hypotheses using a Chi-squared test of association

Chi-squared Test of Association

• Thus, we have two 2x2 tables:

Observed:

	Dis	sease	
Risk Factor	Yes	No	Total
Exposed	$a = O_{11}$	$b = O_{01}$	a+b
Non-exposed	$c = O_{10}$	$d = O_{00}$	c+d
Total	a+c	b+d	\overline{n}

Expected:

	Dise		
Risk Factor	Yes	No	Total
Exposed	(a+b)(a+c)/n	(a+b)(b+d)/n	a+b
	$= E_{11}$	$= E_{01}$	
Non-exposed	(c+d)(a+c)/n	(c+d)(b+d)/n	c+d
	$= E_{10}$	$= E_{00}$	
Total	a+c	b+d	\overline{n}

Example

Devel	oned		1
DCVCI	opcu	\sim	u

		•	
Vitamin C	Yes	No	Total
Yes	17	122	139
No	31	109	140
Total	48	231	279

Calculate Expected Counts

$$E_{11} = 23.9$$

$$E_{11} = 23.9$$
 $E_{01} = 115.1$

$$E_{10} = 24.1$$

$$E_{10} = 24.1$$
 $E_{00} = 115.9$

How to construct the test?

Chi-Squared Test

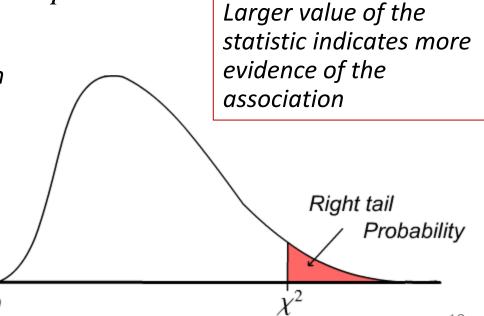
Hypothesis: H_0 : there <u>is no</u> association H_1 : there <u>is</u> an association

Test statistic:

$$X^{2} = \sum \frac{(observed - expected)^{2}}{expected}$$

Follows the chi-square distribution with degrees of freedom:

(# rows - 1)x(# cols - 1)



Fisher's Exact Test

- The p-value for a chi-squared test of association is approximate; the approximation does poorly in "small" samples
- \bullet A common rule-of-thumb for "small" is if any of the (expected) cell counts is ≤ 5
- In small samples, we do not use a large-sample approximation, but compute the exact *p*-value
- This approach was first proposed by R.A. Fisher, hence the name Fisher's Exact Test

Motivating Example

• Example: A retrospective case/control study among men aged 50-54 who died over a 1-month period; their dietary habits were ascertained from a close relative

	Salt I		
Disease	High	Low	Total
CVD	5	30	35
No CVD	2	23	25
Total	7	53	60

How to construct the test?

Fisher's Exact Test

	Dise		
Risk Factor	Yes	No	Total
Exposed	?	?	a+b
${\sf Non\text{-}exposed}$?	?	c+d
Total	a+c	b+d	n

- Fisher's Exact Test fixes the row and column totals (margins)
- Once we know one of the cells (a,b,c,d), we know them all because the margins are fixed
 - (a) Determine every possible table that would still lead to the same margins as those observed
 - (b) Compute the probability that each of the tables would be observed
 - (c) p-value = sum of probabilities corresponding to our observed table and any tables more extreme than our observed table

Analysis Using R

```
> riskratio.wald(M,rev = c("both"))
$data
       cold
vitc
          N Y Total
        109 31
                 140
 Ν
  Υ
        122 17
               139
  Total 231 48
                 279
$measure
    risk ratio with 95% C.I.
vitc estimate
                   lower
                             upper
                                NA
   N 1.0000000
                      NA
   Y 0.5523323 0.3209178 0.9506203
$p.value
    two-sided
vitc midp.exact fisher.exact chi.square
                          NA
   Ν
             NA
                                      NA
   Y 0.02951602
                  0.03849249 0.02827186
```

```
> oddsratio.wald(M,rev = c("both"))
$data
       cold
          N Y Total
vitc
        109 31
                 140
        122 17
               139
  Total 231 48
                 279
$measure
    odds ratio with 95% C.I.
vitc estimate
                   lower
                             upper
   N 1.0000000
                      NΑ
                                NΑ
  Y 0.4899524 0.2569419 0.9342709
$p.value
    two-sided
vitc midp.exact fisher.exact chi.square
             NΑ
                          NΑ
                  0.03849249 0.02827186
   Y 0.02951602
```

Paired Data

McNemar's test for paired binary data

Use exact McNemar's test for small data

R code

- https://www.statology.org/mcnemars-test-r/
- https://yuzar-blog.netlify.app/posts/2022-02-20mcnemar/