

# Comparing Proportions in two-by-two tables:

This section presents measures for comparing groups on binary responses.

The following table is from a report on the relationship between aspirin use and heart attacks.

A five-year randomized study was conducted to test "whether regular intake of aspirin reduces mortality from heart attack disease". Every other day, the participants took either one aspirin tablet or a placebo. The study was "blinded" (the participants did not know which type of pill they were taking).

Group	Heart attack		Total
	Yes	No	
Placebo	189	10,845	11,034
Aspirin	104	10,933	11,037

Now we want to compare the success probabilities (Heart attack disease) of two groups. For this purpose we have three methods:

## (1) Difference of Proportion:

We treat groups as independent binomial sample.

In which  $n_1 = 11,034$ , 189 participants who took placebo suffered from Heart attack

and their proportion is  $\frac{189}{11,034} = 0.0171$ .

And  $n_2 = 11,037$ , 104 participants who took aspirin are suffered from heart attack. and their proportion is  $\frac{104}{11,037} = 0.0094$ .

The sample difference of proportions is

$$\begin{aligned} & P_1 - P_2 \\ &= 0.0171 - 0.0094 \\ &= 0.0077 \end{aligned}$$

(taking aspirin appears to result in a diminished risk of heart attack).

The difference  $P_1 - P_2$  always fall between  $-1$  and  $+1$ .

we compare it with zero.

If  $P_1 - P_2 > 0$  ( $P_1$  has more risk) or  $P_2$  has less risk

If  $P_1 - P_2 < 0$  ( $P_1$  has less risk) or  $P_2$  has more risk

If  $P_1 - P_2 = 0$  ( $P_1$  &  $P_2$  has risk effect)

⇓

It means the response is independent of the group classification.

$$C.I = (P_1 - P_2) \pm Z_{\alpha/2} (SE) < \pi_1 - \pi_2 < (P_1 - P_2) \pm Z_{\alpha/2} (SE)$$

$$S.E = \sqrt{\frac{P_1(1-P_1)}{n_1} + \frac{P_2(1-P_2)}{n_2}}$$

## (2) Relative Risk / Risk Ratio

Another way to compare proportions is

$$\text{Relative Risk} = \frac{P_1}{P_2} \quad \text{where } P_2 \neq 0$$

compare it by 1

If  $\frac{P_1}{P_2} > 1$  ( $P_1$  has more risk)

$\frac{P_1}{P_2} < 1$  ( $P_2$  has more risk)

$\frac{P_1}{P_2} = 1$  ( $P_1$  &  $P_2$  has equal risk)

$$\text{Relative Risk} = \frac{0.0171}{0.0094} = 1.82$$

(The sample proportion of heart attack cases was 82% higher for the group taking placebo.)

The sample difference of proportions of 0.0077 makes it seem as if the two groups differ by a trivial amount. By contrast, the relative risk shows that the difference may have important implications. {Using the difference of proportions alone to compare two groups can be misleading when the proportions are both close to zero.}

The sampling distribution of the sample relative risk is highly skewed unless the sample sizes are quite large. Because of this, its C.I. is rather complex.

A large sample confidence interval for the log of the relative risk is

$$\log(P_1/P_2) \pm Z_{\alpha/2} \sqrt{\frac{1-P_1}{n_1 P_1} + \frac{1-P_2}{n_2 P_2}}$$

Antilogs of the endpoints yield an interval for the true relative risk.

The 95% C.I. to the true relative risk  $\pi_1/\pi_2$  is (1.43, 2.30).

### 3-Relative Odds / Odds Ratio

odd+1

$$\text{Odds} = \frac{P_1}{1-P_1} \text{ or } \frac{P_2}{1-P_2}$$

Odds of having heart attack for those persons who take placebo

$$\begin{aligned} \text{Odds}_1 &= \frac{P_1}{1-P_1} = \frac{0.0171}{1-0.0171} = \frac{0.0171}{0.9829} \\ &= 0.0174 \text{ (1.74\% yes outcome)} \\ &\text{for every 100 no} \end{aligned}$$

odds of heart attack for those persons who take aspirin.

$$\begin{aligned} \text{Odds}_2 &= \frac{P_2}{1-P_2} = \frac{0.0094}{1-0.0094} = \frac{0.0094}{0.9906} \\ &= 0.0094 \text{ (0.94\% yes outcome)} \\ &\text{for every 100 no} \\ &\text{outcomes).} \end{aligned}$$

$$\begin{aligned} \text{Odds Ratio} &= \frac{\text{Odds}_1}{\text{Odds}_2} = \frac{P_1/1-P_1}{P_2/1-P_2} = \frac{\pi_{11}/\pi_{12}}{\pi_{21}/\pi_{22}} \\ &= \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}} \\ &= \frac{0.0174}{0.0094} = 1.85 \end{aligned}$$

The odds value of placebo is 1.85 times or 85% higher than the aspirin.

- odd Ratio = 1 (odds<sub>1</sub> is same as odds<sub>2</sub>)
- odd Ratio < 1 (odds<sub>2</sub> is higher than odds<sub>1</sub>)
- odd Ratio > 1 (odds<sub>1</sub> is higher than odds<sub>2</sub>)

Sample odds ratio

$$\hat{O} = \frac{P_{11}/1 - P_{12}}{P_{21}/1 - P_{22}} = \frac{n_{11}/n_{12}}{n_{21}/n_{22}} = \frac{n_{11} n_{22}}{n_{12} n_{21}}$$

→ The odds ratio does not change value when the table oriented reverses so that the rows become the columns and columns become the rows. The same value occurs when we treat the columns as the response variable and rows as the explanatory variable or vice versa.

Thus it is unnecessary to identify one classification as a response variable in order to estimate odds ratio. By contrast, the relative risk requires this.

So the odds ratio can be defined using joint probabilities as

$$\text{Odds Ratio} = \frac{\pi_{11}/\pi_{12}}{\pi_{21}/\pi_{22}} = \frac{\pi_{11} \pi_{22}}{\pi_{12} \pi_{21}}$$

The odds ratio is also called the Cross-product ratio because it equals the ratio of the products  $\pi_{11} \pi_{22}$  and  $\pi_{12} \pi_{21}$  of cell probabilities from diagonally opposite cells.

→ The sample odds ratio equal 0 or ∞ if  $n_{ij} = 0$  for any  $n_{ij}$ .  
**Inference for odds ratios and log odds ratios.**

Unless the sample size is extremely large, the sampling distribution of the odds ratio is highly skewed.

Because of this skewness, statistical inference for the odds ratio uses an alternative but equivalent measure  $\Rightarrow$  Its natural logarithm (log odd ratio).

log odds ratio = 0 (independence)

The log odds ratio is symmetric about zero, in the sense that reversing rows or reversing columns changes its sign.

The sample log odd ratio ( $\log \hat{\theta}$ ) has a less skewed sampling distribution that is bell-shaped. Its approximating normal distribution has a mean of  $\log \theta$  and a standard error of

$$S.E = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

The S.E. decreases as the cell counts increase.

Because sampling distribution is closer to normality for  $\log(\hat{\theta})$  than  $\hat{\theta}$  so it is better to construct confidence intervals for  $\log \theta$ . Transform back (that is, take antilogs, using the exponential function) to form a confidence interval for  $\theta$ .

A large-sample confidence interval for  $\log \hat{O}$  is

$$CI(\log \hat{O}) = \log \hat{O} \pm Z_{\alpha/2} (S.E)$$

Exponentiating endpoints of this confidence interval yields one for  $O$ .

## Relationship b/w odds ratios and Relative Risk:

Relative Risk = 1.83 means that  $P_1$  is 1.83 times <sup>higher than</sup>  $P_2$  or  $P_1$  is 83% higher than  $P_2$ .

When the odds ratio = 1.83 means that odds value  $\frac{P_1}{1-P_1}$  is 1.83 times <sup>higher than</sup> the odds value  $\frac{P_2}{1-P_2}$  or  $\frac{P_1}{1-P_1}$  is 83% higher than  $\frac{P_2}{1-P_2}$ .

$$\text{Odds ratio} = \frac{P_1/1-P_1}{P_2/1-P_2} = \text{Relative Risk} \times \left( \frac{1-P_2}{1-P_1} \right)$$

When  $P_1$  and  $P_2$  are both close to zero, the fraction in the last term of this expression equals approximately 1. The odds ratio and Relative Risk then take similar values.

This relationship between the odds ratio and the relative risk is useful for some data sets direct



estimation of the relative risk is not possible, yet one can estimate the odd ratio and use it to approximate the relative risk.

## Independence:

Two variables are said to be statistically independent if the population conditional distribution of  $Y$  are identical at each level of  $X$ .

e.g.: The probability of having car is equal for both females and males.

When both variables are response variables, we can describe their relationship using their joint distribution or the conditional distribution of  $Y$  given  $X$ , or the conditional distribution of  $X$  given  $Y$ . Statistically independence is equivalently the property that all joint

probabilities equal the product of their marginal probabilities.

$$\pi_{ij} = \pi_{i+} \cdot \pi_{+j} \quad \text{for } i = 1, \dots, I \\ j = 1, \dots, J$$

Thus the probability that  $X$  falls in row  $i$  and  $Y$  falls in column  $j$  is the product of the probability that  $X$  falls in row  $i$  with the probability the  $Y$  falls in column  $j$ .