



Odds Ratio

Sayantee Jana

Odds Ratio

$$OR = Odds_1 / Odds_2 = \frac{p_1/(1-p_1)}{p_2/(1-p_2)}$$

Alternative measure

- **Relative Risk(RR)** : $RR = p_1 / p_2$

Why use OR?

- Easy to interpret
- Popular since in most practical situations we have binary data
- **invariant under study design**
- Good **mathematical properties**



Example :

Cross-Classification of Smoking by Lung Cancer

Smoking status	Lung Cancer	
	Cases	Controls
Y	688	650
N	21	59
Total	709	709

- Odds of LC for smokers = $688/21 = \Omega_1$,(say)
- Odds of No LC for smokers = $650/59 = \Omega_2$,(say)
- **OR = $\Omega_1 / \Omega_2 = 2.97 \approx 3$**
- $\text{s.e.}(\log(\text{OR})) = \sqrt{\frac{1}{688} + \frac{1}{650} + \frac{1}{21} + \frac{1}{59}} = 3.847$
- P value = $\Pr(\log(\text{OR})/\text{s.e.}(\text{OR}) > \text{obs value} \mid H_0) = 0$

Interpretations

- $OR=1 \Rightarrow$ disease and exposure status are **independent**
- $OR>1 \Rightarrow$ smokers are more likely to have LC than non-smokers i.e. **Positive association** between exposure and disease
- $OR<1 \Rightarrow$ non-smokers are more likely to have LC than smokers i.e. **negative association** between exposure and disease

Software

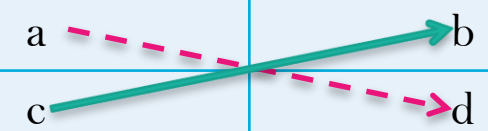
R :Command: `summary(glm(Y ~ X , data = ,
family=binomial(link="logit")))`



Different ways of
calculating OR

From contingency tables

	Y	
X	Y=0	Y=1
X=0	a	b
X=1	c	d



- $H_0 : OR=1$
- $OR = (a*d)/(b*c)$

$$s.e.(\log(OR)) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

- $Z_{\text{teststat}} = \log(OR) / s.e.(OR)$

Logistic Regression Model: Gender vs Status

Let $\pi_0 = P(J | F)$, $\pi_1 = P(J | M)$

$$\log\{\pi/(1-\pi)\} = \text{intercept} + \text{gender} * \beta$$

joining status= 1, if joined
joining status=0, if not joined

gender = 1 , F
= 0 , M

$$\log \text{ of odds}(\pi_0) = \text{intercept} + \beta$$

$$\log \text{ of odds}(\pi_1) = \text{intercept}$$

$$\begin{aligned} \log \text{ odds ratio comparing F to M} \\ &= \log \text{ odds for F} - \log \text{ odds for M} \\ &= (\text{intercept} + \beta.1) - (\text{intercept} + \beta.0) \end{aligned}$$

$$H_0: \log(\text{OR}) = 0 \approx \beta = 0$$

Maximum likelihood equation

- $m = \text{no. of } F$
- $n = \text{no of } M$

$$l = \prod_{1}^m \pi_0 \prod_{1}^n \pi_1$$

Why Logistic Regression instead of 2×2 tables?

- multiple categories.
- multiple covariates
- Dependent variable is always a binary outcome.
- Independent variables may be categorical or quantitative.

