

- Contingency Tables
- More Exercises
- Discuss Quizzes/Answer Questions

11.0 Contingency Tables

This formula is due to the Rev. Thomas Bayes, 1763.

$$P[A_i|B] = \frac{P[B|A_i] * P[A_i]}{\sum_{i=1}^n P[B|A_i] * P[A_i]}$$

$P[A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_n] = 1$. Then

Let A_1, \dots, A_n be mutually exclusive and suppose that

Recall the simple version of Bayes' Rule::

11.1 More Exercises

Suppose Karl Rove takes a lie detector test about outing CIA agent Valerie Plame. And suppose a lie detector is 92% accurate on truth tellers, but only 85% accurate on liars.

Also suppose that the leak had to come from one of five people, all of whom are equally suspect.

If Karl Rove denies leaking but the lie detector signals, then what is the probability that he was responsible for the leak?

Thus the guilty signal from the lie detector only increases the probability that Rove was the leak from $\frac{1}{5}$ to .726.

$$= .726$$

$$= (.85) * (.2) / [(.85) * (.2) + (1 - .92) * (.8)]$$

$$\frac{P[\text{signal} \mid \text{lie}] * P[\text{lie}]}{P[\text{signal} \mid \text{lie}] * P[\text{lie}] + P[\text{signal} \mid \text{truth}] * P[\text{truth}]}$$

$$P[\text{lie} \mid \text{signal}] =$$

Set $B = \{\text{signal}\}$, $A_1 = \{\text{lie}\}$, $A_2 = \{\text{truth}\}$.

Male	Female	Math	10	20	English	20	10	History	15	15
------	--------	------	----	----	---------	----	----	---------	----	----

A **contingency table** shows counts for two categorical variables. For example, you might classify a sample of people by gender and major, or by race and marital status, or by diet and health.

11.2 Contingency Tables

The question is whether membership in the aspirin group is independent of whether one is in the heart attack group, or whether aspirin affects (lowers) the heart attack rate.

Heart Attack	No Heart Attack	aspirin	placebo
104	10933	189	10845

The Physician's Health Study (1989) was a randomized, controlled, double-blind experiment. Half of 22071 men over 40 took an aspirin every other day, while half took a placebo.

One way to assess whether or not treatment is independent of heart attack is to make a mosaic plot. In a mosaic plot, the size of the tiles is proportional to the count in the cell. If it looks like a window, one has perfect independence.

But we know that you are unlikely to get exactly equal numbers of heads for both groups. By chance, one or the other will have a slight excess of heart attacks.

For example, if you toss a coin to determine who gets a heart attack and who does not, then aspirin use is clearly unrelated to heart attack.

But no data set is really going to give you a perfect windowpane. Even if the two categories are independent, random chance will cause some variation.

- Odds Ratio
 - Relative Risk
- For now, we shall describe two ways of measuring the how far from independence a particular 2×2 contingency table is.
- Determining whether or not an observed contingency table shows significant difference from independence is something we shall study later.
- We need some way to measure how unlikely the observed data are, if in fact aspirin use has nothing to do with heart attacks.
- For the study of aspirin and heart attacks, we can determine whether there is a significant difference between the observed and expected frequencies in each cell of the 2×2 table.
- For example, consider the following table:
- | | Yes | No | Total |
|-------|-----|-----|-------|
| Yes | 100 | 100 | 200 |
| No | 100 | 100 | 200 |
| Total | 200 | 200 | 400 |
- The total number of observations is 400. The total number of 'Yes' responses is 200. The total number of 'No' responses is 200. The total number of 'Yes' responses who used aspirin is 100. The total number of 'Yes' responses who did not use aspirin is 100. The total number of 'No' responses who used aspirin is 100. The total number of 'No' responses who did not use aspirin is 100.
- Now, let's calculate the expected frequencies for each cell of the table. The expected frequency for the top-left cell ('Yes', 'Yes') is calculated as follows:
- $$\text{Expected Frequency} = \frac{\text{Row Total} \times \text{Column Total}}{\text{Grand Total}}$$
- $$= \frac{200 \times 200}{400}$$
- $$= 100$$
- Similarly, the expected frequency for the top-right cell ('Yes', 'No') is 100, the expected frequency for the bottom-left cell ('No', 'Yes') is 100, and the expected frequency for the bottom-right cell ('No', 'No') is 100.
- Now, let's calculate the Chi-squared statistic for this table:
- $$\chi^2 = \sum \frac{(O - E)^2}{E}$$
- $$= \frac{(100 - 100)^2}{100} + \frac{(100 - 100)^2}{100} + \frac{(100 - 100)^2}{100} + \frac{(100 - 100)^2}{100}$$
- $$= 0$$
- Since the Chi-squared statistic is 0, we fail to reject the null hypothesis of independence. This means that there is no significant difference between the observed and expected frequencies in the table.

This is just the ratio of the proportion in Category 1 in the first row to the proportion in Category 1 in the second row.

$$RB = \frac{A/(A+B)}{C/(C+D)}$$

Relative Risk is

	Category 1		Category 2		Type 2
	A	B	C	D	
Type 1					

To define these measures, consider the following 2×2 contingency table:

is $1/3$.

The odds are just the ratio of the proportion in group to the proportion in the other, where the sum of the proportions must add to 1. Thus odds of "2 to 1" mean that the first outcome is twice as likely as the second, or the chance of the first outcome is $2/3$, while the chance of the second

the to odds of being in Category 1 in the second row.

This is the ratio of the odds of being in the Category 1 in the first row

$$RR = \frac{C/D}{A/B}$$

The Odds Ratio is

and the treatment of interest.

It is important to get the “A” cell right—it should be the bad outcome

$$OR = \frac{189/10845}{104/10933} = .546.$$

$$RR = \frac{189/(189 + 10845)}{104/(104 + 10933)} = .55$$

Placebo	189	10845
aspirin	104	10933
Heart Attack	No Heart Attack	

In the case of aspirin data:

If the mosaic plot is a perfect windowpane, then both the Relative Risk and the Odds Ratio equal to 1. But if either is much less than 1, that is strong evidence of dependence (i.e., aspirin helps). If the first row people are more likely to be in the first column than second row people, then both Relative Risk and the Odds Ratio are greater than 1. If the first row people are less likely to be in the first column than second row people, then both Relative Risk and the Odds Ratio are less than 1. Thus one can tell whether aspirin hurts, helps, or is independent of heart attack.

As in our example, it often happens that both Relative Risk and the Odds Ratio have similar, but not identical, values.

Statisticians slightly prefer Relative Risk as a measure of dependence, since it is interpretable has the ratio of two probabilities.

For historical reasons, the Odds Ratio is often used in medicine and gambling. But it is slightly more difficult to interpret, at least for me.