

Week Four

Last Week

- Bayesian Inference for Binomial & Multinomial Distributions
- Bayesian Inference vs. Maximum Likelihood Estimation
- Contingency Table Primer Video
 - Contingency Table Overview: multiple categorical variables
 - joint, marginal, and conditional probabilities

This Week: Contingency Tables

Today:

- Activity
 - Comparing Proportions: Relative Risk & Odds Ratios
 - Chi-Squared Tests for Independence
- Thursday: Lab

Next Week: More Contingency Tables

Comparing Proportions

Consider two binary variables. As an example, participants are given ice cream from either sweet peaks or genuine and asked whether it was delicious (5 stars). Note: this could be displayed in a 2 X 2 contingency table.

	Yes	NO	
SP	$\pi_{sp,y}$	$\pi_{sp,n}$	π_{sp}
Gen	$\pi_{g,y}$	$\pi_{g,n}$	π_g

We may be interested in comparing the proportion of respondents that rated ice cream as delicious given the ice cream shop that made it. Note these are conditional probabilities:

$$\pi_{y|sp} = \frac{\pi_{sp,y}}{\pi_{sp}}.$$

Here we can directly compare $\pi_{y|sp}$ and $\pi_{y|g}$. Recall for binomial data ($Y \sim \text{Binomial}(n, p)$) that:

- $E[Y] = np$
- $\text{Var}[Y] = np(1 - p)$

We can use the MLE estimates of $\pi_{y|sp}$, $\pi_{y|g}$, and $\pi_{y|sp} - \pi_{y|g}$

- MLE of $\pi_{y|sp}$: $p_{y|sp} = \frac{n_{sp,y}}{n_{sp}}$, where is n
- MLE of $\pi_{y|g}$: $p_{y|g} = \frac{n_{g,y}}{n_g}$
- MLE of $\pi_{y|sp} - \pi_{y|g}$: $p_{y|sp} - p_{y|g}$

We can use a large sample approximation to construct a confidence interval.

- The SE for $p_{y|sp} - p_{y|g}$ is $\sqrt{\frac{p_{y|sp}(1-p_{y|sp})}{n_{sp}} + \frac{p_{y|g}(1-p_{y|g})}{n_g}}$

$$p_{y|sp} - p_{y|g} \pm z_{\alpha/2}(SE)$$

Example: Construct an uncertainty interval for the difference in proportions when,

- $n_{sp,y} = 80$

- $n_{g,y} = 65$

- $n_{sp} = n_g = 100$

```
n_spy <- 80
n_gy <- 65
n <- 100
p_sp <- n_spy / n
p_gy <- n_gy / n
diff <- p_sp - p_gy
SE <- sqrt(p_sp * (1 - p_sp) / n + p_gy * (1 - p_gy) / n)
```

A 95% interval can be calculated as (0.028, 0.272)

Sometimes we are interested in other comparisons of binary proportions, consider the ratio of the success probabilities $\frac{\pi_{sp}}{\pi_g}$.

- The ratio of probabilities is referred to as relative risk and is fairly common in medical settings. Consider the two sets of probabilities: 0.410 and 0.401 versus 0.010 and 0.001. Calculate the difference and relative risk
- difference: the differences for the first pair is 0.009, which is the same (0.009) as the second pair.
- relative risk: the relative risks are 1.02 and 10 respectively, which are very different numbers.
- Is it a good thing or bad thing that differences and relative risks can have such contrasting implications?

Another comparison, which we will see in much more detail later involves odds.

- $\text{odds} = \pi/(1 - \pi)$,

so consider the following probabilities and odds:

- $\pi = .5$, odds = 1
- $\pi = \frac{2}{3}$, odds = 2, which implies that the odds of a success are twice as likely as a failure.
- $\pi = .75$, odds = 3
- $\pi = .8$, odds = 4

- While odds are useful for looking at a single event, we are often interested in comparing two binary events. In addition to differences and relative risk, we can also consider the odds ratio.

- The odds ratio is defined as $\frac{\text{odds}_1}{\text{odds}_2} = \frac{\pi_1(1-\pi_1)}{\pi_2(1-\pi_2)}$

- We will return to odds ratios in the context of logistic regression in coming weeks.

χ^2 test for independence

Recall the table created in the video lectures

	NO	EB	
SP	76	47	123
Gen	35	42	77
	111	89	200

Our question was whether circadian rhythm and ice cream preference were independent.

Generically we can write this table as

	1	2	
1	π_{11}	π_{12}	π_{1+}
2	π_{21}	π_{22}	π_{2+}
	π_{+1}	π_{+2}	

where the π values that include a + are marginal values.

- Then independence can be stated as $H_0 : \pi_{ij} = \pi_{i+}\pi_{+j} \forall i, j$. Describe this in words, we are saying that if the variables are independent the joint probability can be calculated directly from the marginal values. In contrast, if the joint probability cannot be derived directly from the marginal values, then there is dependence as the likelihood of the outcome of one variable depends on the other.

To test this hypothesis (of independence), we can use the Pearson χ^2 statistic,

$$\chi_{df}^2 = \sum \frac{(n_{ij} - \mu_{ij})^2}{\mu_{ij}}$$

where the degrees of freedom is $(I - 1) \times (J - 1)$.

The μ_{ij} values can be calculated as $n * \pi_{i+} * \pi_{+j}$.

	NO	EB
SP	68.265	54.735
Gen	42.735	34.265

for comparison, the observed data counts were

	NO	EB
SP	76	47
Gen	35	42

For reference, our test statistic is 5.116, which would result in a p-value of 0.024.

```
chisq.test(as.matrix(output_table[1:2,2:3]))
```

Pearson's Chi-squared test with Yates' continuity correction

```
data: as.matrix(output_table[1:2, 2:3])
X-squared = 4.4757, df = 1, p-value = 0.03438
```

```
chisq.test(as.matrix(output_table[1:2,2:3]), simulate.p.value = T)
```

Pearson's Chi-squared test with simulated p-value (based on 2000 replicates)

```
data: as.matrix(output_table[1:2, 2:3])
X-squared = 5.1157, df = NA, p-value = 0.02149
```

Recall, we know the true values

	1	2	
1	.4	.2	.6
2	.2	.2	.4
	.6	.4	

and $\pi_{11} = .4 \neq .36 = \pi_{1+}\pi_{+1}$

Finally, let's think about this problem in the context of estimation. We can directly estimate the probabilities (joint, marginal, or conditional) associated with this data.

Assuming we use a uniform Dirichlet prior, estimate posterior means and intervals for the four joint probabilities.

```
dirichlet_samples <- as_tibble(rdirichlet(10000, c(1 + 76, 1 + 47, 1 + 35, 1 + 42)),
                              .name_repair = c('minimal'))
colMeans(dirichlet_samples)
```

```
0.3764865 0.2351203 0.1769809 0.2114123
```

```
apply(dirichlet_samples, 2, quantile, prob=c(.025, .975))
```

```
2.5% 0.3120672 0.1795763 0.1268578 0.1585053
97.5% 0.4430419 0.2980894 0.2321389 0.2719695
```