

Up until now  $\rightarrow$  asymptotic inference

2/8/22

\* Requires large samples

$\rightarrow$  CLT  $\rightarrow$  approximate CIs & p-values using normal approximation

$R$ : prop.test

- One sample & 1 binary variable  $\rightarrow$  exact binomial inference (p-value)  
 $R$ : binom.test

Exact inference for 2 proportions (2 binary variables):

① Fisher's Exact test

② Randomization (simulation-based) tests \* Not in book

## Fisher's Exact Test

Example: Penguins -

Scope of Inference  $\rightarrow$  ① Cause-&-effect?



Randomized experiment

Observational Study

Yes

No

② To whom can we generalize?

Representative sample of what population?

- Only generalize to penguins similar to those in the sample

Do totals fixed  $\rightarrow$  "binomial" sampling

Data:

	Survive	Not	
Metal	3	7	10
Not metal	6	4	10
	9	11	20

Consider count in 1<sup>st</sup> cell:  $n_{11} = 3$

$H_0$ : No association between survival status - whether the penguin had a metal band.

$$\pi_1 = P(\text{survive} | \text{metal band})$$
$$\pi_2 = P(\text{survive} | \text{no metal band})$$

$$H_0: \pi_1 - \pi_2 = 0$$

$$H_0: \frac{\pi_1}{\pi_2} = 1$$

$$\text{Odds ratio} = \theta = \frac{\pi_1 / (1 - \pi_1)}{\pi_2 / (1 - \pi_2)}$$

$$H_0: \frac{\pi_1 / (1 - \pi_1)}{\pi_2 / (1 - \pi_2)} = 1$$

$$H_0: \theta = 1$$

$$H_a: \theta < 1$$

Under  $H_0 \Rightarrow$  Model  $n_{11}$  with a hypergeometric distribution!

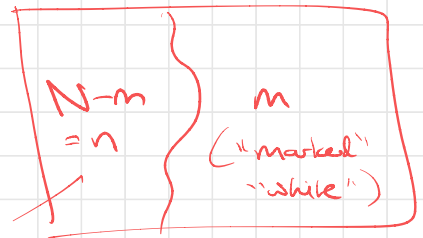
Parameters:

$m$  = # marked in pop.

$n$  = # unmarked in pop.

$k$  = sample size

Scenario: Finite population of  $N$  items.



"unmarked" "black"  $N$  items

- Sample  $k$  items without replacement.

$$X \sim \text{Bin}\left(k, \frac{m}{n+m}\right)$$

$X$  = # of marked items in the sample

	Success	Failure	
Group 1	$n_{11}$	$n_{12}$	$n_{1+}$
Group 2	$n_{21}$	$n_{22}$	$n_{2+}$
	$n_{+1}$	$n_{+2}$	$n$

$n$  Sample  $n_{1+}$   
from  
population  
of  
 $n$  items

Assuming  $H_0$ : No association -

$n_{+1}$  marked  
 $n_{+2}$  unmarked

$X$  = # of marked items  
that end up in our sample.

$$P(X = x) = \frac{\text{\# of samples where } x \text{ marked } \& \text{ } k-x \text{ unmarked}}{\text{total \# of possible samples}}$$

$$= \frac{\binom{m}{x} \binom{n}{k-x}}{\binom{m+n}{k}} = \frac{\binom{n_{+1}}{x} \binom{n_{+2}}{k-x}}{\binom{n}{n_{+1}}}$$

$R$ : dhyper( $x, m, n, k$ )

dhyper( $n_{11}, n_{+1}, n_{+2}, n_{+1}$ )

$$p\text{-value} = P(X \leq 3)$$

$R$ : phyper(3, 9, 11, 10)

col. 1  
total

col. 2  
total

row 1 total

Randomization Test - Simulate the random assignment process 1000's of times assuming  $H_0$ .

p-value = proportion of those simulations which resulted in the observed data or something more extreme

Ex:

	Survive	Not	
Metal	3	7	10
No metal	6	4	10
	9	11	20

Under  $H_0$   $\hookrightarrow$  9 penguins would survive regardless of group assignment

Simulation: 9 survivors }  $\rightarrow$  10  
 11 non-survivors }  $\rightarrow$  10  
 R.A.

How to simulate w/ cards?

Total # of cards? 20 - 1 for each penguin

blue (survivors) 9  
 red (non-survivors) 11

Shuffle

10 metal 10 non-metal