#### Lecture 21

Ciprian M Crainiceanu

Table of contents

Outline

Fisher's exa test

The hyperged metric distribution

test in practice

Monte Carlo

### Lecture 21

### Ciprian M Crainiceanu

Department of Biostatistics Johns Hopkins Bloomberg School of Public Health Johns Hopkins University

December 5, 2013

### Table of contents

## Table of contents

Outline

test

The hyperged metric distribution

test in practic

Monte Carl

- 1 Table of contents
- 2 Outline
- 3 Fisher's exact test
- 4 The hypergeometric distribution
- 5 Fisher's exact test in practice
- 6 Monte Carlo

Table of contents

Outline

Fisher's exa test

The hyperged metric distribution

Fisher's exact test in practice

Monte Carl

- 1 Introduce Fisher's exact test
- 2 Illustrate Monte Carlo version of test

The hypergeo metric distribution

test in practice

Monte Carlo

# Fisher's exact test

- Fisher's exact test is "exact" because it guarantees the  $\alpha$  rate, regardless of the sample size
- Example, chemical toxicant and 10 mice

	Tumor	None	Total
Treated	4	1	5
Control	2	3	5
Total	6	4	

- $p_1$  = prob of a tumor for the treated mice
- $p_2$  = prob of a tumor for the untreated mice

Fisher's exact test in practice

Monte Carl

### Continued

scone 12 18 1 3 8 8 18 N 15 short

- $H_0: p_1 = p_2 \neq p$
- Can't use Z or  $\chi^2$  because SS is small
- Don't have a specific value for p

Fisher's exact

Monte Carl

### Fisher's exact test

- Under the null hypothesis every permutation is equally likely
- observed data

Treatment : T T T T T C C C C C Tumor : T T T T N T T N N N

• permuted back the connection W. fix mangin

Treatment: TCCTCTTCTC

JTJC

Tumor : N T T N N T T T N T 67 4A

• Fisher's exact test uses this null distribution to test the

• Fisher's exact test uses this null distribution to test the hypothesis that  $p_1 = p_2$ 

The hypergeometric distribution

Fisher's exact test in practice

Monte Carlo

# Hyper-geometric distribution

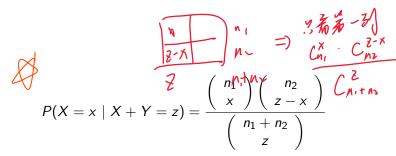
- X number of tumors for the treated
- Y number of tumors for the controls
- $H_0: p_1 = p_2 = p$
- Under  $H_0$ 
  - $X \sim \text{Binom}(n_1, p)$
  - *Y* ∼ Binom(*n*<sub>2</sub>, *p*)
  - $X + Y \sim \mathsf{Binom}(n_1 + n_2, p)$

Fisher's exa

The hypergeometric distribution

test in practice

Monte Carlo



This is the hypergeometric pmf

metric distribution

test in practice

Monte Carlo

$$P(X = x) = \binom{n_1}{x} p^x (1 - p)^{n_1 - x}$$

$$P(Y = z - x) = \binom{n_2}{z - x} p^{z - x} (1 - p)^{n_2 - z + x}$$

$$P(X + Y = z) = \binom{n_1 + n_2}{z} p^z (1 - p)^{n_1 + n_2 - z}$$

The hyperge

The hypergeometric distribution

test in practice

Monte Carlo

$$P(X = x \mid X + Y = z) = \frac{P(X = x, X + Y = z)}{P(X + Y = z)}$$

$$= \frac{P(X = x, Y = z - x)}{P(X + Y = z)}$$

$$= \frac{P(X = x)P(Y = z - x)}{P(X + Y = z)}$$

Plug in and finish off yourselves

Table of contents

Outille

The hypergeo

Fisher's exact test in practice

Monte Carlo

- More tumors under the treated than the controls.
- Calculate an exact P-value
- Use the conditional distribution = hypergeometric
- Fixes both the row and the column totals
- Yields the same test regardless of whether the rows or columns are fixed
- Hypergeometric distribution is the same as the permutation distribution given before

The hypergeo metric distribution

Fisher's exact test in practice

Monte Carlo

# Tables supporting $H_a$

Consider  $H_a: p_1 > p_2$ 

"me sided"

- P-value requires tables as extreme or more extreme (under  $H_a$ ) than the one observed
- Recall we are fixing the row and column totals
- Observed table

Table 
$$1 = \begin{array}{c|c} 4 & 1 & 5 \\ \hline 2 & 3 & 5 \\ \hline 6 & 4 & \end{array}$$

The Mangah

More extreme tables in favor of the alternative

Table 2 = 
$$\begin{bmatrix} 5 & 0 & 5 \\ 1 & 4 & 5 \end{bmatrix}$$

Table of

Outline

test

The hyperged metric distribution

Fisher's exact test in practice

Monte Carlo

P(Table 1) = 
$$P(X = 4|X + Y = 6)$$
  
=  $\begin{pmatrix} 5\\4 \end{pmatrix} \begin{pmatrix} 5\\2 \end{pmatrix}$   
=  $\begin{pmatrix} 10\\6 \end{pmatrix}$ 

P(Table 2) = 
$$P(X = 5|X + Y = 6)$$
  
=  $\begin{pmatrix} 5 \\ 5 \end{pmatrix} \begin{pmatrix} 5 \\ 1 \end{pmatrix}$   
=  $\begin{pmatrix} 10 \\ 6 \end{pmatrix}$ 

P-value = 0.238 + 0.024 = 0.262

Fisher's exact test in practice

R code

dat <- matrix(c(4, 1, 2, 3), 2)  $\begin{pmatrix} \zeta & 1 \\ 2 & 3 \end{pmatrix}$  fisher.test(dat, alternative = "greater")  $\chi = \frac{2}{\zeta}$ 

Ho: P1=P2 Fisher's Exact Test for Count Data

data: dat.

p-value = 0.2619

alt hypoth: true odds ratio is greater than 1

-----output-----

95 percent confidence interval:

0.3152217 Tnf sample estimates:

odds ratio

4.918388

The hypergeo

Fisher's exact test in practice

Monte Carl

- Two sided p-value = 2×one sided P-value (There are other methods which we will not discuss)
- P-values are usually large for small n
- Doesn't distinguish between rows or columns
- The common value of p under the null hypothesis is called a nuisance parameter
- Conditioning on the total number of successes, X + Y, eliminates the nuisance parameter, p
- Fisher's exact test guarantees the type I error rate
- Exact unconditional P-value

$$\sup_{p} P(X/n_1 > Y/n_2; p)$$

Fisher's exact test

The hypergeo metric distribution

Fisher's exact test in practice

Monte Carlo

Observed table X = 4

Treatment: T T T T T C C C C C
Tumor: T T T T N T T N N N

Permute the second row

Treatment: TTTTTCCCCCC
Tumor: TNTNTTNNTT

- Simulated table X = 3
- Do over and over
- Calculate the proportion of tables for which the simulated
   X 4
- This proportion is a Monte Carlo estimate for Fisher's exact P-value