

## Statistical Inference for Two Proportions.

Our parameters of interest:

$p_i, i=1,2 \dots$  the (sub)population proportion for the (sub)population  $i=1,2$

e.g.,  $p_1$  corresponds to the subpopulation who get the placebo;  
 $p_2$  corresponds to the subpopulation who get the treatment.

Sample sizes are denoted by  $n_1$  and  $n_2$ .

Assume that the two samples are independent.

For large  $n_i, i=1,2$ , we know the approximate dist'n of:

- the count r.v.s:

$$X_i \approx \text{Normal}(\text{mean} = n_i \cdot p_i, \text{sd} = \sqrt{n_i \cdot p_i(1-p_i)}) \quad i=1,2$$

and

- the sample proportion r.v.s:

$$\hat{P}_i = \frac{X_i}{n_i} \approx \text{Normal}(\text{mean} = p_i, \text{sd} = \sqrt{\frac{p_i(1-p_i)}{n_i}}) \quad i=1,2$$

We base our statistical inference for  $p_1 - p_2$  on  $\hat{P}_1 - \hat{P}_2$ .

$$\hat{P}_1 - \hat{P}_2 \approx \text{Normal}(\text{mean} = p_1 - p_2, \text{sd} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}})$$

## Confidence Interval

### C... confidence level

$$\text{pt. estimate} \pm \text{margin of error}$$

$$z^* \cdot \text{stderror}$$

$$z^* = \Phi^{-1}\left(\frac{1+C}{2}\right) \\ = qnorm\left(\frac{1+C}{2}\right)$$

$$\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$z^* \cdot \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

w/  $\hat{p}_1, \hat{p}_2$

$\hat{p}_1 - \hat{p}_2$  are the observed sample proportions

## Hypothesis Testing.

$$H_0: p_1 = p_2$$

vs.

$$H_a: \begin{cases} p_1 < p_2 \\ p_1 \neq p_2 \\ p_1 > p_2 \end{cases}$$

Test statistic :  $\hat{P}_1 - \hat{P}_2$

Under the null hypothesis :  $p_1 = p_2 = p$

$$\begin{aligned} \hat{P}_1 - \hat{P}_2 &\approx \text{Normal}(\text{mean} = 0, \text{sd} = \sqrt{\frac{p(1-p)}{n_1} + \frac{p(1-p)}{n_2}}) = \\ &= \sqrt{p(1-p) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \end{aligned}$$

$$\frac{\hat{P}_1 - \hat{P}_2}{\sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \approx N(0,1) \text{ under the null hypothesis}$$

Let  $\hat{p}_1$  and  $\hat{p}_2$  be the observed sample proportions.

Q: What's the form of the observed z-statistic?

$$z = \frac{\hat{P}_1 - \hat{P}_2}{\sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

w/ the estimate  $\hat{p}$  based on all the available data.

We pool the two samples to estimate  $p$ , i.e.,

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

w/  $x_i, i=1, 2 \dots$  # of successes  
in sample i

Note:

$$\hat{p} = \frac{n_1 \cdot \hat{p}_1 + n_2 \cdot \hat{p}_2}{n_1 + n_2} = \frac{n_1}{n_1 + n_2} \cdot \hat{p}_1 + \frac{n_2}{n_1 + n_2} \cdot \hat{p}_2$$

To calculate the p-value, we find:

$$\begin{cases} H_0: p_1 < p_2 : & p\text{-value: } P[Z < z] \\ H_a: p_1 \neq p_2 : & p\text{-value: } P[Z < -|z|] + P[Z > |z|] \\ H_a: p_1 > p_2 : & p\text{-value: } P[Z > z] \end{cases}$$

If a significance level  $\alpha$  is given, then ...

if p-value  $\leq \alpha$ , we REJECT THE NULL.

if p-value  $> \alpha$ , we FAIL TO REJECT THE NULL.