2-Sample Categorical Data

BIOSTAT 201A Fall 2025

Discussion 7 – November 14, 2025

Cindy J. Pang

		Proportions	
(₀ =	P1 = P2	(=) Pi-P2 = 0	

$$H_1 = p_1 \neq p_2$$
 $\Leftrightarrow p_1 - p_2 \neq 0$

"success"

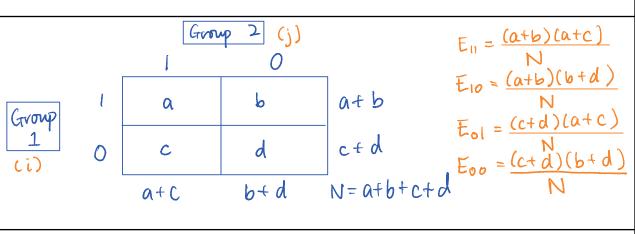
H₀: Group 1 and Group 2 are independent
$$H_1$$
: Group 1 and 2 are associated.

All expected counts are ≥ 5 ($E_1^* \geq 5$)

(2) Rule of Five satisfied for both proportions s.t.
$$n_1p_1(l-p_1) \ge 5$$
 and $n_2p_2(l-p_2) \ge 5$

Binary covaniate

Total



"failure"

$$Z = \frac{\hat{\rho}_{1} - \hat{\rho}_{2}}{\sqrt{\hat{\rho}_{p} (1 - \hat{\rho}_{p}) (\frac{1}{n_{1}} + \frac{1}{n_{2}})}} \text{ where } \hat{\rho}_{p} = \frac{n_{1} \hat{\rho}_{1} + n_{2} \hat{\rho}_{2}}{n_{1} + n_{2}} = \frac{x_{1} + x_{2}}{n_{1} + n_{2}}$$

$$= \frac{(O_{1}) - E_{1}}{E_{1}}^{2} \text{ where } E_{1} = \frac{(row total in i)(col total in j)}{N}$$

$$= \frac{(a - E_{1})^{2}}{E_{1}} + \frac{(b - E_{10})^{2}}{E_{10}} + \frac{(d - E_{00})^{2}}{E_{00}}$$

Hypotheses

Assumptions

Data

$$= \frac{(a - C_{11})^{2} + (b - E_{10})}{E_{10}} + \frac{(c - E_{01})^{2} + (d - E_{00})^{2}}{E_{00}}$$

$$= \frac{(a - C_{11})^{2} + (b - E_{10})}{E_{10}} + \frac{(a - E_{00})^{2}}{E_{00}} + \frac{(a - E_{00})^{2}}{E_{00}}$$

$$= \frac{(a - C_{11})^{2} + (b - E_{10})}{E_{10}} + \frac{(a - E_{00})^{2}}{E_{00}} + \frac{(a - E_{00})^{2}}{E_{00}} + \frac{(a - E_{00})^{2}}{E_{00}}$$

$$= \frac{(a - C_{11})^{2} + (b - E_{10})}{E_{10}} + \frac{(a - E_{00})^{2}}{E_{00}} + \frac{(a -$$

$$= \frac{\hat{p}_{1} - \hat{p}_{2}}{\sqrt{\hat{p}_{p} (1 - \hat{p}_{p}) (\frac{1}{n_{1}} + \frac{1}{n_{2}})}} \text{ where } \hat{p}_{p} = \frac{n_{1} \hat{p}_{1} + n_{2} \hat{p}_{2}}{n_{1} + n_{2}} = \frac{x_{1} + x_{2}}{n_{1} + n_{2}}$$

A hypothesis suggests that breast cancer risk increases with a longer interval between the onset of menstruation and a woman's first childbirth, making age at first birth a potential risk factor. To test this, an international study examined women from hospitals in the United States, Greece, Yugoslavia, Brazil, Taiwan, and Japan. Among women who had given birth, 21.2% of those with breast cancer (683 of 3,220) and 14.6% of those without breast cancer (1,498 of 10,245) had their first child at age 30 or older. The question is whether this difference reflects a true association or simply occurred by chance.

(a) Test this hypothesis using a 2-Sample Test for Difference in Population Proportions

$$H_{6}^{2} P_{BC}, A_{9}e^{30} = P_{BC}, A_{9}e^{<30} \quad \text{vs.} \quad H_{1}^{2} P_{BC}, A_{9}e^{30} \neq P_{BC}, A_{9}e^{<30}$$

$$P_{1} = \frac{A_{1}}{N_{1}} = \frac{683}{3220} = 21.2\% = .212$$

$$P_{2} = \frac{A_{1}}{N_{1}} = \frac{683}{3220} = 21.2\% = .212$$

$$P_{2} = \frac{A_{1}}{N_{2}} = \frac{1.498}{10.245} = 14.6\% = .146$$

$$P_{2} = \frac{A_{2}}{N_{2}} = \frac{1.498}{10.245} = 14.6\% = .146$$

$$P_{3} = \frac{A_{1}}{N_{1}} = \frac{A_{1}}{N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{4} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{2} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{4} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{3} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{4} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{2.181}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + N_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + A_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + A_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

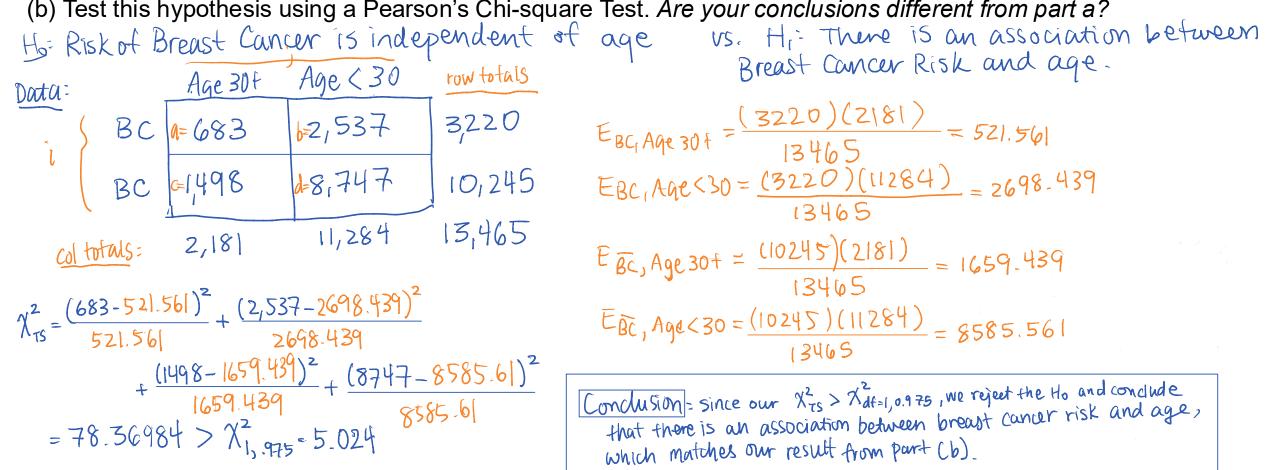
$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + A_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{N_{1} + A_{2}} = \frac{A_{1} + A_{2}}{15.465} = 0.1619755$$

$$P_{5} = \frac{A_{1} + A_{2}}{15.465} = 0.$$

A hypothesis suggests that breast cancer risk increases with a longer interval between the onset of menstruation and a woman's first childbirth, making age at first birth a potential risk factor. To test this, an international study examined women from hospitals in the United States, Greece, Yugoslavia, Brazil, Taiwan, and Japan. Among women who had given birth, 21.2% of those with breast cancer (683 of 3,220) and 14.6% of those without breast cancer (1,498 of 10,245) had their first child at age 30 or older. The question is whether this difference reflects a true association or simply occurred by chance.

(b) Test this hypothesis using a Pearson's Chi-square Test. Are your conclusions different from part a?



Breast Cancer Risk and age.

$$E_{BC_1}Aqe_{30f} = \frac{(3220)(2181)}{13465} = 521.561$$
 $E_{BC_1}Aqe_{30f} = \frac{(3220)(11284)}{13465} = 2698.439$
 $E_{BC_2}Aqe_{30f} = \frac{(10245)(2181)}{13465} = 1659.439$
 $E_{BC_3}Aqe_{30f} = \frac{(10245)(11284)}{13465} = 8585.561$
 $E_{BC_4}Aqe_{30f} = \frac{(10245)(11284)}{13465} = 8585.561$

Conclusion: since our X2 > Xdf=1,0.975, we reject the Ho and conclude

which matches our result from part (b)

that there is an association between breast cancer risk and age,