

Assignment 4

This assignment will provide experience with confidence intervals and hypothesis testing.

1. A Finnish study compared stroke in adults age 50-66 with a history of back pain versus no back pain. Of $n_1=1212$ with back pain, 180 (14.9%) had stroke in a 5 year follow up period. Of the $n_2=576$ with no back pain, 58 (10.1%) had stroke in the same 5 year period.

The standard error for a single proportion P is given by $SE(P) = \sqrt{P(1-P)/n}$ where n is the sample size.

1a. What is the difference in the two stroke proportions?

1b. Analogous with the calculation for the difference in two independent means, what is the standard error for the proportion in each group and the standard error for the difference in two proportions (calculate & see notes).

1c. Form a 95% confidence interval for the true population difference in the two proportions.

1d. What would be the conventional null value (value under the null hypothesis) of the difference between two proportions for the usual type of hypothesis testing? Does the interval computed above include or exclude this value?

1e. Compute the observed Z statistic for testing the null hypothesis above and give the corresponding one or two sided p value (label).

Under the null hypothesis that the true $\pi_1=\pi_2=\pi$, one may use the common $P = (n_1P_1 + n_2P_2)/(n_1+n_2)$ to compute the standard error for the difference in proportions (SE diff) by substituting this value for P_1 and P_2 . In Excel, the Gaussian distribution percentile for Z is given by the function `=normsdist(Z)`.

1f. Give the corresponding Odds ratio (OR) for stroke in those with back pain vs no back pain. (This is the OR actually observed).

1g. Give a p value for testing the odds ratio above versus the null hypothesis that the true underlying population OR=1 or briefly explain why this cannot be done (computations not required).

1h. Will the corresponding 95% confidence interval for the true odds ratio contain 1.0? (Computations not required).

Question 2 requires use of a table in Section V of the notes (you determine which table) or software such as GPower (free software)

2. You are planning a study to compare pregnancy rates in women aged 30-35 treated with a new fertility drug relative to an old drug. This will be a randomized clinical trial. In previous studies using the old drug, after one cycle of treatment, 25% of women became pregnant. You wish to confirm at least a 20% increase (ie to at least 45% pregnant) under the new drug after one cycle.

2a. What sample size is needed per group to achieve 80% power using the usual two sided $p < 0.05$ significance criterion?

2b. Does the required sample size increase or decrease if 90% power is required?

2c. Does the sample size increase or decrease if $\alpha=0.10$ is used instead of $\alpha=0.05$ as the two sided significance criterion.

2d. Does the sample size increase or decrease if the pregnancy rate is 30% under the old drug but is still 45% for the new drug? Does this correspond to greater or less patient homogeneity within each group?

2e. If the sample size is 30 per group, what is the smallest increase from 25% that can be confirmed with 80% power using $\alpha = 0.05$.

3 The sample correlation coefficient, r , measures the association between two continuous variables measured on the same subjects (This will be studied later in class). When there is no correlation, the true population correlation value, ρ , equals **zero** (null value, correlation = $\rho=0$).

In a Swedish study, the sample correlation (r) between mothers height and their infants birth weight in a sample of $n=30$ women was $r=0.40$ with a corresponding 95% confidence interval of (0.04, 0.67).

3a. Indicate which are true or false

- i. The corresponding two sided p value is > 0.05
- ii. The corresponding two sided p value is < 0.05
- iii. The corresponding two sided p value is > 0.10
- iv. The corresponding two sided p value is < 0.10
- v. It is not possible to know the p value without more information

The quantity r^2 , the square of the correlation, represents how much of the variation in birth weight that is accounted for by mothers height.

3b. What is the upper bound on how much of the variation in birthweight is accounted for by mothers height (a proxy for mothers “size”).

4 Give the correct test (A, B, C D, E or F) for computing a p value corresponding to each statistic comparison

Tests: A: t test, B: Chi square test/Fisher exact test, C: log rank test D: Wilcoxon rank sum test E: ANOVA F test F: Kruskal-Wallis test

4a: Comparing two or more proportions - test is ____

4b: Comparing two means – test is ____

4c: Comparing two or more survival curves – test is ____

4d: Comparing three or more means (overall test) – test is ____

4e: Comparing two or more medians (non normal continuous data) – test is ____