

**Data Analysis**  
Spring Semester, 2023  
February 2, 2023  
Lecture 4

Syllabus Quiz deadline is extended to February 2, 11:59 pm.

**Feb 23: Examination One: Chapters 3, 4, 5, 6, and 7**

Two major points:

In a randomized experiment, a difference in response between the group given an experimental treatment and the group given a control treatment is **CAUSED** by the experimental treatment or is a chance event.

In an observational study, correlation shows association, not necessarily causation.

***Chapter 4, Probability and Probability Distributions***

*Conditional Probability and Bayes' Theorem*

Definition of conditional probability of A given B:  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ ,  $P(B) > 0$ .

The sets A, B, and C are a collection of cover sets if the sample space S is such that  $S = A \cup B \cup C$ , where A, B, and C are disjoint (mutually exclusive or incompatible)

The law of total probability:  $P(E) = P(E|A)P(A) + P(E|B)P(B) + P(E|C)P(C)$ .

Bayes' Theorem:  $P(A|E) = \frac{P(E|A)P(A)}{P(E)} = \frac{P(E|A)P(A)}{P(E|A)P(A) + P(E|B)P(B) + P(E|C)P(C)}$ .

**Normal Probability Distribution.**

Let the random variable X be normally distributed with expected value  $\mu$  and variance  $\sigma^2$ .

That is,  $X \sim N(\mu, \sigma^2)$ . Any probability calculation about a normal distribution can be transformed to a calculation with a standard normal:

$$P(X \leq a) = P\left(\frac{X - \mu}{\sigma} \leq \frac{a - \mu}{\sigma}\right) = \Phi\left(\frac{a - \mu}{\sigma}\right),$$

where  $\Phi$  is the cumulative distribution function of the standard normal. That is,

$\Phi(z) = \Pr \{Z \leq z\}$ , where  $Z \sim N(0,1)$ .

Key Percentiles of the Standard Normal:

- $P(Z \leq -2.576) = \Phi(-2.576) = 0.005$
- $P(Z \leq -2.326) = \Phi(-2.326) = 0.01$
- $P(Z \leq -1.960) = \Phi(-1.960) = 0.025$
- $P(Z \leq -1.645) = \Phi(-1.645) = 0.05$
- $P(Z \leq -1.282) = \Phi(-1.282) = 0.10$
- $P(Z \leq -0.6745) = \Phi(-0.6745) = 0.25$

Make sure that you understand how to use Table 1 of the statistical tables in the back of your text and in the tables that you will use for your examinations.

Definition of Expected Value of a Discrete Random Variable:

$$E(X) = \sum_{APV} xP(X = x),$$

where *APV* means to sum over all possible values of the discrete random variable.

Expectation is a linear operator.

Definition of variance of the random variable  $X$ :  $var(X) = E((X - EX)^2)$

Important identity:

$$var(X) = E((X - EX)^2) = E(X^2) - (EX)^2$$

The proof starts with

$$\begin{aligned} E[(X - EX)^2] &= E[X^2 - 2X(EX) + (EX)^2] \\ &= E[X^2] - E[2X(EX)] + E[(EX)^2] \\ &= E[X^2] - 2(EX)E[X] + (EX)^2 \\ &= E[X^2] - 2(EX)^2 + (EX)^2 \\ &= E(X^2) - (EX)^2. \end{aligned}$$

Example Problem:

The random variables  $W_1$  and  $W_2$  are a random sample of 2 drawn from the random variable  $W$  which has expected value  $\mu_W$  and standard deviation  $\sigma_W > 0$ . Find  $E(W_1 - W_2)$  and  $E((W_1 - W_2)^2)$ .

Solution:

$$E(W_1 - W_2) = E(W_1) - E(W_2) = \mu_W - \mu_W = 0;$$

These steps follow from the linear operator property of expectation.

To find  $E((W_1 - W_2)^2)$ :

$$\begin{aligned} E((W_1 - W_2)^2) &= E[W_1^2 - 2W_1W_2 + W_2^2] \\ &= E[W_1^2] - E[2W_1W_2] + E[W_2^2]. \end{aligned}$$

First, since  $var(W) = E[W^2] - [E(W)]^2$ ,

$$E[W^2] = var(W) + [E(W)]^2 = \sigma_W^2 + \mu_W^2.$$

Second,  $E[W_1W_2] = E(W_1)E(W_2)$ , since  $W_1$  and  $W_2$  are independent.

Then,  $E[W_1^2] - E[2W_1W_2] + E[W_2^2] = \sigma_W^2 + \mu_W^2 - 2E(W_1)E(W_2) + \sigma_W^2 + \mu_W^2$

Combining,

$$\begin{aligned} E((W_1 - W_2)^2) &= 2\sigma_W^2 + 2\mu_W^2 - 2\mu_W^2 \\ &= 2\sigma_W^2. \end{aligned}$$

In a random sample of size 2,  $\sigma_W^2 = \frac{E((W_1 - W_2)^2)}{2}$ . An unbiased estimate of  $\sigma_W^2$  is  $\frac{(W_1 - W_2)^2}{2}$ .

## **Chapter Five**

### **Inferences about Population Central Values**

Distribution of Sample Mean:

- Let  $Y_1, Y_2, \dots, Y_n$  be a random sample of size  $n$  from  $Y$  which has the distribution  $N(\mu, \sigma^2)$ .

**THEN** the distribution of  $\bar{Y}_n = \frac{\sum_{i=1}^n Y_i}{n}$ , the sample mean, is  $N(\mu, \sigma^2/n)$ .

- Central Limit Theorem (CLT): When  $Y_1, Y_2, \dots, Y_n$  is a random sample of size  $n$  from  $Y$  which has expected value  $\mu$  and variance  $\sigma^2 < \infty$ , then the

distribution of  $\bar{Y}_n = \frac{\sum_{i=1}^n Y_i}{n}$  is asymptotically  $N(\mu, \sigma^2/n)$ . For a random sample of size  $n$ , it is always true that  $E(\bar{Y}_n) = \mu$  and  $\text{var}(\bar{Y}_n) = \frac{\sigma^2}{n}$ . The CLT allows probability calculations that increase in accuracy as the sample size increases.

Testing a Statistical Hypothesis (Variance Known)

- Null hypothesis:  $H_0 : E(Y) = \mu_0$
- Alternative Hypothesis:  $H_1 : E(Y) \neq \mu_0$
- Level of significance  $\alpha$
- Type I error: reject a null hypothesis that is true.
- The probability of a Type I error is  $\alpha$ . Formally,  $\Pr_0\{\text{Reject } H_0\} = \alpha$ .
- Test statistic. Null distribution is the distribution of the test statistic under the null hypothesis.
- Type II error: accept a null hypothesis that is false.
- Typically,  $\alpha$  is set to a small number (0.05 or 0.01), and  $n$  is chosen so that  $\beta = \Pr_1\{\text{Accept } H_0\}$ , where  $\beta$  is dependent on a setting of the alternative hypothesis, is small. Alternative distribution: distribution of the test statistic under the setting of the alternative hypothesis.

Testing a Statistical Hypothesis (Variance Unknown).

- We cannot put  $\bar{Y}_n$  in standard score form:  $Z = \frac{\bar{Y}_n - \mu_0}{\sigma / \sqrt{n}}$  because we do not know  $\sigma$ .

- Instead, we use an estimate of  $\sigma^2$ ,  $\hat{\sigma}^2 = S^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}{n-1}$ , which has  $n-1$  degrees of freedom.

- We put  $\bar{Y}_n$  in studentized standard score form:  $T_{n-1} = \frac{\bar{Y}_n - \mu_0}{\hat{\sigma} / \sqrt{n}}$ .

- Student showed that the percentiles from the standard normal (here 1.645, 1.960, and 2.576) had to be stretched. The amount of stretching is

determined by the number of degrees of freedom in  $\hat{\sigma}^2 = S^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}{n-1}$ ;

here,  $n-1$  degrees of freedom. These values are tabulated and will be given to you in your examinations.

- Example problem: Chapter 5 Study Guide, Problem 4:

A research team took a sample of 8 observations from the random variable  $Y$ , which had a normal distribution  $N(\mu, \sigma^2)$ . They observed  $\bar{y}_8 = 43.2$ , where  $\bar{y}_8$  is the average of the eight sampled observations and  $s^2 = 517.5$  is the observed value of the unbiased estimate of  $\sigma^2$ , based on the sample values. Test the null hypothesis that  $H_0 : E(Y) = 50$  against the alternative  $H_1 : E(Y) \neq 50$  at the 0.10, 0.05, and 0.01 levels of significance.

- The degrees of freedom is  $n-1 = 8-1 = 7$ .
- The studentized test statistic is  $t_7 = \frac{43.2 - 50}{\sqrt{(517.5/8)}} = \frac{-6.8}{\sqrt{64.7}} = \frac{-6.8}{8.04} = -0.845$
- Find the student t stretches for 1.645, 1.960, 2.576. They are 1.895, 2.365, 3.499.
- Make your decision. Here, it is to accept at the 0.10 level of significance (and of course at 0.05 and 0.01 as well) since  $|-0.845| < 1.895$ .

### Confidence Interval, Variance Unknown

- As always, we use the data  $S^2$  to estimate  $\sigma^2$  and stretch our normal percentile (here 2.576) using the degrees of freedom of  $S^2$ . The stretch of 2.576 is 3.499.

- The 99% confidence interval for  $E(Y)$  is  $\bar{y}_n \pm t_{n-1, 2.576} \frac{\hat{\sigma}}{\sqrt{n}}$ .

- In the example problem, the 99% confidence interval for  $E(Y)$  is

$$\bar{y}_8 \pm t_{7, 2.576} \frac{\hat{\sigma}}{\sqrt{n}} = 43.2 \pm 3.499 \sqrt{\frac{517.5}{8}} = 43.2 \pm 3.499(8.04) = 43.2 \pm 28.1.$$

### Paired t-test

The most common application of the paired t-test is a comparison of the post-training score of a participant in a study with the same participant's pre-training score. The idea of the paired t-test is to calculate the difference (here post score-pre score) for each participant.

This data is used in the one-sample t-test of Chapter 5 to test the null hypothesis that the expected post training score is equal to the expected pre training score.

### Chapter 6 Study Guide, Problem 5

A research time wished to estimate the reduction of the density of contaminant in a liquid due to filtering the liquid. They filtered four samples, called A, B, C, and D. Find the 99% confidence interval for the expected reduction in the density of contaminant using the data in the table below:

Sample	Density of Contaminant before Filtering	Density of Contaminant after Filtering	Difference	Deviation of Difference	Deviation Squared
A	132	87	45	-4.75	22.5625
B	205	163	42	-7.75	60.0625
C	81	35	46	-3.75	14.0625
D	423	357	66	16.25	264.0625

Solution: The four differences of  $D$ =before filtering – after filtering are: 132-87=45, 42, 46, and 66. Then  $\bar{d}_4 = \frac{45+42+46+66}{4} = 49.75$  and  $s_D^2 = 120.25$  on 3 degrees of freedom. The 99% confidence interval for the expected reduction in density is

$$49.75 \pm 5.841 \sqrt{\frac{120.25}{4}},$$

which is the interval from 17.7 to 81.8.

REJE

CT



## Chapter Six

### Inferences Comparing Two Population Central Values

*New Problem (just slightly different from Chapter 4, number 7):*

The random variables  $W_1$  and  $W_2$  are independent. Their expected values are  $E(W_1) = \mu_1$  and  $E(W_2) = \mu_2$ . Their variances are  $\text{var}(W_1) = \sigma_1^2 < \infty$ , and  $\text{var}(W_2) = \sigma_2^2 < \infty$ . Find  $E(W_1 - W_2)$  and  $\text{var}(W_1 - W_2)$ .

**Solution:**

First,  $E(W_1 - W_2) = E(W_1) - E(W_2) = \mu_1 - \mu_2$ .

Second,  $\text{var}(W_1 - W_2) = E\{[(W_1 - W_2) - E(W_1 - W_2)]^2\}$

$$\text{var}(W_1 - W_2) = E\{[(W_1 - W_2) - (\mu_1 - \mu_2)]^2\}$$

$$\text{var}(W_1 - W_2) = E\{[(W_1 - \mu_1) - (W_2 - \mu_2)]^2\}$$

$$\text{var}(W_1 - W_2) = E\{[(W_1 - \mu_1)^2 + (W_2 - \mu_2)^2 - 2(W_1 - \mu_1)(W_2 - \mu_2)]\}$$

$$\text{var}(W_1 - W_2) = E[(W_1 - \mu_1)^2] + E[(W_2 - \mu_2)^2] - 2E[(W_1 - \mu_1)(W_2 - \mu_2)]$$

$$\text{var}(W_1 - W_2) = \text{var}(W_1) + \text{var}(W_2) - 2\text{cov}(W_1, W_2)$$

$$\text{var}(W_1 - W_2) = \sigma_1^2 + \sigma_2^2 - 2 \cdot 0$$

$$\text{var}(W_1 - W_2) = \sigma_1^2 + \sigma_2^2$$

### Two Independent Sample Test

Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from the random variable  $X$ , which is  $N(\mu_X, \sigma_X^2)$ . For example,  $X$  could be the response of a participant in a clinical trial to a new medicine. Let  $B_1, B_2, \dots, B_m$  be a random sample of size  $m$  from the random variable  $B$ , which is  $N(\mu_B, \sigma_B^2)$ . Continuing the example,  $B$  could be the response of a participant in a clinical trial to the best available medicine. The two samples are independent of each other. The context of this discussion is that there is random assignment of the participants in the study to the two groups. This has the effect of roughly balancing the two groups with respect to each and every variable other than the treatment assigned. If there is a significant difference between the two groups, that difference is either a chance event or the causal result of a difference in the treatments.

Back to the probability theory of the analysis,  $\bar{X}_n$  is  $N(\mu_X, \frac{\sigma_X^2}{n})$ , and  $\bar{B}_m$  is  $N(\mu_B, \frac{\sigma_B^2}{m})$

The two sample averages are independent. We seek to use this data to test

$H_0 : E(X - B) = 0$  against the alternative hypothesis  $H_1 : E(X - B) \neq 0$  at the  $\alpha$  level of significance. Our test statistic is  $TS = \bar{X}_n - \bar{B}_m$ .

### *Distribution of the Test Statistic*

The distribution of  $TS = \bar{X}_n - \bar{B}_m$  is normal. In the case that the random variable  $X$  is not normal (or the random variable  $B$  is not normal), then the CLT implies that the distribution of  $TS$  is approximately normal. From the problem at the start of class,  $E(TS) = E(\bar{X}_n - \bar{B}_m) = E(\bar{X}_n) - E(\bar{B}_m) = \mu_X - \mu_B$ , and

$$\text{var}(TS) = \text{var}(\bar{X}_n - \bar{B}_m) = \text{var}(\bar{X}_n) + \text{var}(\bar{B}_m) = \frac{\sigma_X^2}{n} + \frac{\sigma_B^2}{m}.$$

### *Null Distribution of the Test Statistic*

Since  $H_0 : E(X - B) = 0$  which is equivalent to  $H_0 : \mu_X = \mu_B$ ,

$$E_0(TS) = E_0(\bar{X}_n - \bar{B}_m) = \mu_X - \mu_B = 0.$$

R. A. Fisher argued that the null hypothesis should be that the random variable  $X$  has exactly the same distribution as the random variable  $B$ . That is, not only does  $\mu_X = \mu_B$  under the null, but also  $\sigma_X^2 = \sigma_B^2 = \sigma^2$ . In fact, any probabilistic parameter of random variable has the same value for  $X$  and  $B$ . Since it is very unusual for a treatment to have a measurable effect, this conception of the null hypothesis is widely used. Using this assumption,  $\text{var}_0(TS) = \text{var}_0(\bar{X}_n - \bar{B}_m) = \frac{\sigma^2}{n} + \frac{\sigma^2}{m} = \sigma^2(\frac{1}{n} + \frac{1}{m})$ .

### *Test when variance known*

As in Chapter 5, we test this null hypothesis by putting  $TS$  in standard score form:

$$Z = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{\sigma^2(\frac{1}{n} + \frac{1}{m})}}.$$

If  $\alpha = 0.10$ , we reject  $H_0$  when  $|Z| \geq 1.645$ . If  $\alpha = 0.05$ , we reject  $H_0$  when  $|Z| \geq 1.960$ .

If  $\alpha = 0.01$ , we reject  $H_0$  when  $|Z| \geq 2.576$ .

### *Test when variances unknown but equal*

Just as in Chapter 5, we use an estimate of  $\sigma^2$  and stretch the critical values an amount determined by the degrees of freedom of our estimate. There are many

estimates of  $\sigma^2$ . For example,  $S_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n-1}$  and  $S_B^2 = \frac{\sum_{i=1}^m (B_i - \bar{B}_m)^2}{m-1}$  are unbiased

estimates of  $\sigma^2$ , with  $n-1$  and  $m-1$  degrees of freedom respectively. That is,

$$E(S_X^2) = E(S_B^2) = \sigma^2. \text{ We use both estimates. Let } S_P^2 = \frac{(n-1)S_X^2 + (m-1)S_B^2}{n+m-2}.$$

This estimator has  $n+m-2$  degrees of freedom. Then our studentized statistic is

$$T_{n+m-2} = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{S_P^2 \left( \frac{1}{n} + \frac{1}{m} \right)}}.$$

If  $\alpha = 0.10$ , we reject  $H_0$  when  $|T_{n+m-2}| \geq t_{1.645, n+m-2}$ . Similarly, if  $\alpha = 0.05$ , we reject  $H_0$  when  $|T_{n+m-2}| \geq t_{1.960, n+m-2}$ . If  $\alpha = 0.01$ , we reject  $H_0$  when  $|T_{n+m-2}| \geq t_{2.576, n+m-2}$ .

*Example Problem: from examination 1, Fall 2016, 2B*

A research team took a random sample of 3 observations from a normally distributed random variable  $Y$  and observed that  $\bar{y}_3 = 37.4$  and  $s_y^2 = 42.6$ , where  $\bar{y}_3$  was the average of the three observations sampled from  $Y$  and  $s_y^2$  was the unbiased estimate of  $\text{var}(Y)$  (i.e., the divisor in the variance was  $n-1$ ). A second research team took a random sample of 5 observations from a normally distributed random variable  $X$  and observed that  $\bar{x}_5 = 50.6$  and  $s_x^2 = 48.1$ , where  $\bar{x}_5$  was the average of the five observations sampled from  $X$  and  $s_x^2$  was the unbiased estimate of  $\text{var}(X)$  (i.e., the divisor in the variance was  $n-1$ ). Test the null hypothesis  $H_0 : E(X) = E(Y)$  against the alternative  $H_1 : E(X) \neq E(Y)$  at the 0.10, 0.05, and 0.01 levels of significance using the pooled variance t-test. This problem is worth 40 points.

Solution:  $s_P^2 = \frac{(5-1)48.1 + (3-1)42.6}{5+3-2} = 46.27$  on 6 degrees of freedom. The t stretches of 1.645, 1.960, and 2.576 for 6 degrees of freedom are 1.943, 2.447, and 3.707 respectively. The test statistic is  $t_{n+m-2} = \frac{37.4 - 50.6 - 0}{\sqrt{46.27 \left( \frac{1}{5} + \frac{1}{3} \right)}} = -2.637$ . Reject the null hypothesis at the 0.10 and 0.05. Accept the null hypothesis at the 0.01 level.

*Alternate Problem:*

A research team took a random sample of 3 observations from a normally distributed random variable  $Y$  and observed that  $\bar{y}_3 = 37.4$  and  $s_y^2 = 42.6$ , where  $\bar{y}_3$  was the average of the three observations sampled from  $Y$  and  $s_y^2$  was the unbiased estimate of  $\text{var}(Y)$  (i.e., the divisor in the variance was  $n-1$ ). A second research team took a random sample of 5 observations from a normally distributed random

variable  $X$  and observed that  $\bar{x}_5 = 50.6$  and  $s_x^2 = 48.1$ , where  $\bar{x}_5$  was the average of the five observations sampled from  $X$  and  $s_x^2$  was the unbiased estimate of  $\text{var}(X)$  (i.e., the divisor in the variance was  $n-1$ ). Find the 99% confidence interval for  $E(X) - E(Y)$ .

Solution: The template for the 99% confidence interval for  $E(X) - E(Y)$  is

$\bar{x}_n - \bar{y}_m \pm t_{2.576, n+m-2} s_P \sqrt{\frac{1}{n} + \frac{1}{m}}$ . That is, the 99% confidence interval for  $E(X) - E(Y)$  is

$$50.6 - 37.4 \pm 3.707 \sqrt{46.27} \sqrt{\frac{1}{5} + \frac{1}{3}} = 13.2 \pm 3.707 \cdot 6.802 \cdot 0.7303 = 13.2 \pm 18.4. \text{ That is, the}$$

99% confidence interval for  $E(X) - E(Y)$  is the interval between -5.2 and 31.6.

Since 0 is in the 99% confidence interval, we should accept  $H_0 : E(X) = E(Y)$  against  $H_1 : E(X) \neq E(Y)$  at the 0.01 level of significance.

### *Test when variances unknown and unequal*

This procedure is called the unequal variance t-test or unequal variance confidence interval. As before, let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from the random variable  $X$ , which is  $N(\mu_X, \sigma_X^2)$ . Let  $B_1, B_2, \dots, B_m$  be a random sample of size  $m$  from the random variable  $B$ , which is  $N(\mu_B, \sigma_B^2)$ . Here,  $\sigma_X^2 \neq \sigma_B^2$ . Then,

$$\text{var}(\bar{X}_n - \bar{B}_m) = \frac{\sigma_X^2}{n} + \frac{\sigma_B^2}{m}.$$

The standard score form of the test statistic is then

$$Z = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_B^2}{m}}}.$$

The studentized standard score form of test statistic would be

$$T_\gamma = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{\frac{S_X^2}{n} + \frac{S_B^2}{m}}}.$$

The problem is that the probability theory calculations that are the basis of the degrees of freedom in the pooled variance t-test are not valid. There is, however, a complicated formula (called Satterthwaite's formula) that produces an approximate degrees of freedom. I will not ask you to calculate this by hand in an examination. Fortunately, the statistical programs that you will use will calculate the unequal variance t-test and unequal variance confidence interval for you. Typically, the equal variance t-test and unequal variance t-tests have essentially equal p-values

when the assumption of equal variance appears reasonable. When the assumption does not appear reasonable, the one should use the unequal variance calculations. My practice is to report the unequal variance results as calculated by a reputable statistics program.

### *Type II Error Rate and Power Calculations*

The definition of the Type II error rate is  $\beta = \Pr\{\text{Accept } H_0\}$ . The power of a statistical test is defined to be  $\text{Power} = 1 - \beta$ . This is the probability that the null hypothesis is correctly rejected. A large Type II error rate indicates a study that is "underpowered." The calculation of  $\beta$  is just a normal probability calculation. The specification of the normal distribution is based on the alternative specified in the problem. Typically, the values of the variances are assumed known.

### *Example Problem: from Chapter 6 Study Guide, Problem 3*

In a clinical trial, 50 patients suffering from an illness will be randomly assigned to one of two groups so that 25 receive an experimental treatment and 25 receive the best available treatment. The random variable  $X$  is the response of a patient to the experimental medicine, and the random variable  $B$  is the response of a patient to the best currently available treatment. The random variables  $X$  and  $B$  are normally distributed with  $\sigma_X = \sigma_B = 500$  under both the null and alternative distributions. The null hypothesis to be tested is that  $E(X) - E(B) = 0$  against the alternative that

单边问题  $E(X) - E(B) > 0$  at the 0.01 level of significance. What is the probability of a Type II error for the test of the null hypothesis when  $E(X) - E(B) = 500$ ?

Solution: The standard score form of the TS is  $Z = \frac{\bar{X}_n - \bar{B}_m - 0}{\sqrt{\sigma^2(\frac{1}{n} + \frac{1}{m})}}$ , and

$H_0 : E(X - B) = 0$  is rejected when  $Z \geq 2.326$ , remembering that the problem asks for a one-sided test at level of significance 0.01. Using the TS directly,

$H_0 : E(X - B) = 0$  is rejected when

方差

$\bar{X}_{25} - \bar{B}_{25} \geq 0 + 2.326 \sqrt{\frac{500^2}{25} + \frac{500^2}{25}} = 0 + 2.326 \cdot 141.42 = 328.95$ . For the alternative specified in the problem  $\bar{X}_{25} - \bar{B}_{25}$  is  $N(500, 141.42^2)$ . Then

$$\beta = \Pr\{\text{Accept } H_0\} = \Pr\{\bar{X}_{25} - \bar{B}_{25} < 328.95\} = \Pr\left\{\frac{\bar{X}_{25} - \bar{B}_{25} - E_1(\bar{X}_{25} - \bar{B}_{25})}{\sigma_1(\bar{X}_{25} - \bar{B}_{25})} < \frac{328.95 - 500}{141.42}\right\}.$$

That is,  $\beta = \Pr\{Z < \frac{328.95 - 500}{141.42} = -1.210\} = \Phi(-1.210) = 0.113$ .

### Sample size for two sample test:

The bad news is that the mathematics of sample size calculations is relatively complex. The good news is that one can solve a wide range of sample size problems once one knows how to solve this one. The argument is essentially the same.

### Problem 6 in Chapter 6 Study Guide

In a clinical trial,  $2J$  patients suffering from an illness will be randomly assigned to one of two groups so that  $J$  will receive an experimental treatment and  $J$  will receive the best available treatment. The random variable  $X$  is the response of a patient to the experimental medicine, and the random variable  $B$  is the response of a patient to the best currently available treatment. The random variables  $X$  and  $B$  are normally distributed. The null hypothesis to be tested is that  $E(X) - E(B) = 0$  against the alternative that  $E(X) - E(B) > 0$  at the  $\alpha$ ,  $\alpha \leq 0.5$ , level of significance. When the null hypothesis is true,  $\text{var}(X) = \text{var}(B) = \sigma_0^2$ . When the alternative hypothesis is true,  $\text{var}(B) = \sigma_0^2$ , but  $\text{var}(X) = \sigma_1^2 > \sigma_0^2$ . What is the number  $J$  in each group that would have to be taken so that the probability of a Type II error for the test of the null hypothesis specified in the common section is  $\beta$ ,  $\beta \leq 0.5$ , when  $E(X) - E(B) = \Delta > 0$ ?

Solution: The test statistic is  $TS = \bar{X}_J - \bar{B}_J$ , and  $TS$  is

$N(E(X) - E(B), \frac{\text{var}(X)}{J} + \frac{\text{var}(B)}{J})$ . The null distribution of  $TS$  is then

$N(0, \frac{\sigma_0^2}{J} + \frac{\sigma_0^2}{J})$ . Hence, we reject  $H_0 : E(X) - E(B) = 0$  against  $H_1 : E(X) - E(B) > 0$

at the  $\alpha$  level of significance when  $TS \geq 0 + |z_\alpha| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_0^2}{J}}$ . When

$E(X) - E(B) = \Delta > 0$  and  $\text{var}(X) = \sigma_1^2 > \sigma_0^2$ ,  $\text{var}(B) = \sigma_0^2$ , the (alternative)

distribution of  $TS$  is then  $N(\Delta, \frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J})$ . Then, the probability of a Type II

error is  $\beta = \Pr\{\text{Accept } H_0\} = \Pr\{TS < 0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}}\}$ . That is,

$$\beta = \Pr\{TS < 0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}}\} = \Pr\{Z = \frac{TS - \Delta}{\sigma_1(\bar{X}_J - \bar{B}_J)} < \frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}}\}. \text{ Since}$$

$\beta = \Pr\{\text{Accept } H_0\} \leq 0.5$ , it is true that  $\beta = \Pr\{Z < -|z_\beta|\}$ . We now have two equations:

$$\beta = \Pr\left\{Z < \frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}}\right\}, \text{ and}$$

$$\beta = \Pr\{Z < -|z_\beta|\}.$$

The problem is to choose  $J$  so that the probability of a Type II error is a specified value. That is, we should choose  $J$  so that the right-hand sides of the

two equations are equal: 
$$\frac{0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta}{\sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}} = -|z_\beta|.$$

We have to solve for  $J$  in the equation above. This reduces to:

$$0 + |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} - \Delta = -|z_\beta| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}}. \quad \text{That is, } |z_\alpha| \sqrt{\frac{2\sigma_0^2}{J}} + |z_\beta| \sqrt{\frac{\sigma_0^2}{J} + \frac{\sigma_1^2}{J}} = \Delta.$$

Next, solve for  $J$  to get  $\sqrt{J} = \frac{|z_\alpha| \sqrt{2\sigma_0^2} + |z_\beta| \sqrt{\sigma_0^2 + \sigma_1^2}}{\Delta}$ . Since  $J$  has to be an integer, we increase  $J$  to the next integer value.

The scenario can be realistic. Suppose that there are two ratio scale random variables. The one with the greater mean typically has a greater variance. For example, the Poisson and the chi-square distributions have this property. This greater variance requires a somewhat greater sample size in study design.

This calculation assumes that there is no attrition of subjects. Typically, study attrition is large. An attrition rate less than 15% at a follow-up three or more years later is a very low attrition rate. Accounting for attrition is its own modeling effort, which can be very difficult to do well.

#### *Problem 4 in Chapter 6 Study Guide*

In a clinical trial,  $2J$  patients suffering from an illness will be randomly assigned to one of two groups so that  $J$  will receive an experimental treatment and  $J$  will receive the best available treatment. The random variable  $X$  is the response of a patient to the experimental medicine, and the random variable

$B$  is the response of a patient to the best currently available treatment. The random variables  $X$  and  $B$  are normally distributed and have  $\sigma_X = \sigma_B = 500$  under both the null and alternative distributions. The null hypothesis to be tested is that  $E(X) - E(B) = 0$  against the alternative that  $E(X) - E(B) > 0$  at the 0.005 level of significance. What is the number  $J$  in each group that would have to be taken so that the probability of a Type II error for the test of the null hypothesis specified in the common section is 0.01 when  $E(X) - E(B) = 250$ ?

Solution: For this specification,  $\alpha = .005$ , so that  $|z_\alpha| = 2.576$ . Also,  $\beta = .01$ , so that  $|z_\beta| = 2.326$ . With regard to variances,  $\sigma_0^2 = \sigma_1^2 = 500^2$ . Finally,  $\Delta = 250$ . Then, the design equation is

$$\sqrt{J} = \frac{2.576\sqrt{2 \cdot 500^2} + 2.326\sqrt{2 \cdot 500^2}}{250} = \frac{3466.24}{250} = 13.865 = \sqrt{192.24}.$$

That is, there should be at least 193 in each group.

The magnitude of the difference in the expected values is 250, which is half of the assumed standard deviation of each group. This is called a one-half standard deviation effect. To detect a difference equal to 1/2 of the standard deviation, researchers need about 200 in each group. This is 400 total observations. One needs about 50 observations per group to detect a one standard deviation effect. That is, fewer observations are needed to detect a larger effect size.