

**I. TYPE I AND TYPE II ERRORS IN HYPOTHESIS TESTING**

**A. The Concepts of Type I and Type II errors**

- As we've been testing hypotheses, we've been using the phrases "reject the null hypothesis" when the calculated statistic falls in the critical region. Or, "fail to reject the null hypothesis" when the calculated statistic does not fall in the critical region.
- Why don't we just say "accept the alternative hypothesis" or "accept the null hypothesis," respectively?
- We've touched on this a bit already, the framework here elaborates on those ideas:
  - (1) we don't know what the true population mean  $\mu$  is equal to (that's why we're using statistical inference to say something about this mean);
  - (2) only one *true* population mean  $\mu$  exists for the sampling distribution; when we run a hypothesis test, we don't know what this true mean is.
- Think about the various ways we can make the right and wrong decisions:

		<b>"State of the World" – the Truth</b>	
		<i>H<sub>0</sub> is true</i>	<i>H<sub>0</sub> is false</i>
<b>Decision after Hypothesis Test</b>	<i>Fail to Reject H<sub>0</sub></i>	Correct decision: probability = $(1-\alpha)$	Type II error: probability = $\beta$
	<i>Reject H<sub>0</sub></i>	Type I error: probability = $\alpha$	Correct decision: probability = $(1-\beta)$ = "power of the test"

- Ideally, the probability of both Type I and Type II errors will be "small."
- However, there is a direct tradeoff between the two types of error *for a given sample size*: the smaller the significance level ( $\alpha$ ), the larger the probability of not rejecting a false null hypothesis ( $\beta$ ).
- So, we have to assess the risks involved in committing either type of error and use that to balance the types of error probabilities.

## B. Calculating Type II Error Probabilities and the Power of the Test<sup>1</sup>

- Key point: the probability of making a Type II error depends on: (1) the sample size, (2) the significance level, and (3) the true value of the parameter
- The following example shows how to compute a Type II error probability. We'll use a one-tail test. In practice, you can calculate power for two-tail hypothesis tests.

Example: The program designers of a new approach for getting low-skilled clients into stable jobs claim that after completing the program, clients earn \$35,000 per year, on average. An advocacy group is skeptical of this claim, and thinks that the program designers are overselling the effectiveness of their program. The advocacy group thinks that the true mean earnings for participants is actually less than \$35,000 per year. The advocacy group plans to perform this test:

$$H_0: \mu_{\text{trainees}} = \$35,000 \text{ per year}$$

$$H_1: \mu_{\text{trainees}} < \$35,000 \text{ per year}$$

They will use a 5% significance level, and a random sample of  $n=3,000$  trainees. The standard deviation of trainee earnings in the sample is \$10,000.

*Q1:* At what sample mean will the null hypothesis be rejected, using  $\alpha = 0.05$  for a one-tail test? i.e., this is the usual hypothesis test that we are used to carrying out:

$$t_{\text{calc}}: \frac{\bar{X} - 35}{10/\sqrt{3,000}} = -1.645 \Leftrightarrow \bar{X} = 35 - 1.645 \left( \frac{10}{\sqrt{3,000}} \right) = 34.7$$

Decision criterion:

If  $\bar{X} \leq \$34.7$  thousand, Reject  $H_0$ . If  $\bar{X} > \$34.7$  thousand, Fail to Reject  $H_0$ .

*Q2:* What is the probability of making a Type II error if the **true mean earnings of trainees in the program** is \$34,850 per year?

In this situation, the TRUE state of the world is that  $\mu_{\text{trainees}} = \$34.85$  thousand, and the standard error is still  $\left( \frac{10}{\sqrt{3,000}} \right) = 0.182574$ .

But remember we set up a decision rule in Q1 that we would reject  $H_0$  if  $\bar{X} \leq 34.7$ .

---

<sup>1</sup> This section draws from Weiss (2002) *Introductory Statistics*, 6<sup>th</sup> ed

We can think about a new, “true” curve of normally distributed  $\bar{X}$  values about their true mean of \$34.85 thousand.

We know that a Type II error occurs if we do not reject Ho (i.e., if  $\bar{X} > \$34.7$  thousand). The probability of that happening is equal to the probability of our drawing a sample with a mean that is greater than 34.7, where that mean is drawn from a population where the true mean is \$34.85 thousand:

$$t_{calc} : \frac{34.7 - 34.85}{0.182574} = -0.82158$$

Because we’re calculating a t-statistic, we’d like to look up -0.82158 in the *t*-table. However, the *t*-table we have doesn’t have this level of detail. Because we have a large enough sample size, to the point where the *t*-distribution approximates the Normal distribution, we can use the Z-table to look up this value. We see that -0.82 corresponds to an area to the left of Z of 0.2061, and an area to the right of Z of  $(1-0.2061) = 0.7939$ .

So if the true mean earnings of trainees in the program is \$34,850 per year, but we are testing the null Ho:  $\mu_{trainees} = \$35,000$  per year (against an alternative of  $\mu_{trainees} < \$35,000$ ), then the probability of making a Type II error is  $\beta = 0.7939$ .

This means that if we were testing a hypothesis that the true mean earnings of trainees in the program is \$35,000 per year, there would a probability of 0.7939, or about a 79% chance, that we would fail to reject this null hypothesis, if in fact the true mean earnings of trainees is \$34,850 per year.

- The probability of *not* making a Type II error is called **power**, or the **power of the test**:

$$\mathbf{Power} = (1 - \text{probability of Type II error}) = \mathbf{1 - \beta}$$

\*\* The closer to 1 the power is, the better the hypothesis test is at detecting a false null hypothesis.

\*\* The closer to zero the power is, the hypothesis test is not very good at detecting a false null.

\*\* In the previous example, the power of the test is **0.2061**. In other words, this test has relatively low power: it is not very good at detecting a false null.

Q3: What is the probability of making a Type II error if the **true mean earnings of trainees in the program** is \$34,250 per year?

In this situation, the TRUE state of the world is that  $\mu_{\text{trainees}} = \$34.25$  thousand, and the standard error is still  $\left(\frac{10}{\sqrt{3,000}}\right) = 0.182574$ .

But remember we set up a decision rule in Q1 that we would reject  $H_0$  if  $\bar{X} \leq 34.7$ .

We can think about a new, “true” curve of normally distributed  $\bar{X}$  values about their true mean of \$34.25 thousand.

We know that a Type II error occurs if we **do not reject  $H_0$**  (i.e., if  $\bar{X} > \$34.7$  thousand). The probability of that happening is equal to the probability of our drawing a sample with a mean that is greater than 34.7, where that mean is drawn from a population where the true mean is \$34.25 thousand:

$$t_{\text{calc}} : \frac{34.7 - 34.25}{0.182574} = 2.46475$$

Remember that because the sample size is large enough, we can look up 2.46 in the Z-table, and see that this corresponds to an area up to Z of 0.9931, and thus an area beyond Z of 0.0069.

So if the true mean earnings of all trainees in the program is \$34,250 per year, the probability of making a Type II error when our null hypothesis is  $H_0: \mu_{\text{trainees}} = \$35,000$  is:  **$\beta = 0.0069$** . The **power of the test** is equal to 0.9931.

If we are testing a hypothesis that the true mean earnings of trainees in the program is \$35,000 per year, there would be less than a 1% chance that we would fail to reject this null hypothesis, if in fact the true mean earnings of trainees is \$34,250 per year.

- The hard thing is that the true value of  $\mu$  is not known: otherwise, we wouldn't be doing hypothesis tests! Often, researchers will construct a table of powers for various possible values of  $\mu$  for a given sample size (in this case,  $n=3,000$ ) e.g.:

IF the True Mean: $\mu =$	Prob. of Type II error: $\beta$	Power: $1 - \beta$
35.25		
35.10		
35.00		
34.90		
34.80		
34.70		
34.60		
34.60		
34.50		
34.40		
34.30		
34.20		

- Power is often a concern in evaluation research, when detecting the effects of a program may be crucial. The relevant point here is that they may be concerned about the “effect size” of a program – how large of an effect can their sample design pick up, given a particular sample size?
- And speaking of sample size....
- Instead of a sample of  $n=3,000$ , what if instead the group sampled  $n=101$  trainees? We can repeat the questions above, now using this smaller sample size:

*Q1\**: At what sample mean will the null hypothesis be rejected, using  $\alpha = 0.05$  for a one-tail test?

$$t_{calc} : \frac{\bar{X} - 35}{\left(\frac{10}{\sqrt{101}}\right)} = -1.660 \Leftrightarrow \bar{X} = 35 - 1.660\left(\frac{10}{\sqrt{101}}\right) = 33.348$$

Decision criterion:

If  $\bar{X} \leq 33.348$ , Reject  $H_0$ .      If  $\bar{X} > 33.348$ , Fail to reject  $H_0$ .

We could now construct a table now using  $n=201$ :

IF the True Mean: $\mu=$	Prob. of Type II error: $\beta$	Power: $1- \beta$
35.25		
35.10		
35.00		
34.90		
34.80		
34.70		
34.60		
34.50		
34.40		
34.30		
34.20		



- Let's take another example: Suppose we're interested in the average family income in the U.S.

First, suppose we draw a random sample of  $n=3,000$  where the sample mean income is \$33,000 and the standard deviation is  $s = \$16,000$ .

- We can construct an interval around this sample mean:  $\bar{X} \pm 1.96 \sigma_{\bar{x}}$   

$$= 33 \pm 1.96 \left( \frac{16}{\sqrt{3000}} \right) = 33 \pm 0.57255$$

[\$32,427, \$33,573]: We can say with 95% confidence that the population mean is between \$32,427 and \$33,573 per year.

- Now, suppose that the true population mean family income is  $\mu_{true} = \$32,750$  and population standard deviation is \$16,000.

**If I were to conduct a hypothesis test and use a null hypothesis of  $\mu_{H_0} = \$33,000$ , what is the probability that I would correctly reject that null hypothesis (i.e.,  $1-\beta$ )? What is the probability that I would incorrectly fail to reject that null hypothesis (i.e.,  $\beta$ )?**

- To calculate these probabilities, we want to map the lower and upper bounds of the confidence interval calculated above onto the distribution of the true mean.

-- Calculate how many standard errors \$32,427 is away from the true mean of \$32,750:

$$\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}} = \frac{32,427 - 32,750}{16/\sqrt{3,000}} = -1.10418 = t \text{ lower}$$

-- Calculate how many standard errors \$33,573 is away from the true mean of \$32,750:

$$\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}} = \frac{33,573 - 32,750}{16/\sqrt{3,000}} = 2.815816 = t \text{ upper}$$

-- Look up these  $t$ -scores in the  $Z$ -table (because at this sample size, the  $t$  approximates the  $Z$ ) to find the area under the curve: The probability of rejecting the null in the left-hand tail of the true curve (centered on  $\mu_{true} = \$32,750$ ) is approximately equal to 0.1357 (i.e., the area to the left of  $Z = -1.10$ ); and the probability of rejecting the null in the right-hand tail of the true curve is approximately equal to 0.0024 (i.e., the area to the right of  $Z = 2.82$ )

-- Thus, the probability of correctly rejecting the null hypothesis of  $\mu_{H_0} = \$33,000$  if the true population mean is  $\mu_{true} = \$32,750 = 0.1357 + 0.0024 = \mathbf{0.1381}$ . The probability of a **Type II** error under these conditions is  $(1-0.1381) = \mathbf{0.8619}$ .