Chapter 4: Foundations for inference

OpenIntro Statistics, 2nd Edition

Variability in estimates

- Application exercise
- Sampling distributions via CLT

2 Confidence intervals

- 3 Hypothesis testing
- Examining the Central Limit Theorem
- Inference for other estimators
- 6 Sample size and power

Statistical vs. practical significance

Young, Underemployed and Optimistic Coming of Age, Slowly, in a Tough Economy

Young adults hit hard by the recession. A plurality of the public (41%) believes young adults, rather than middle-aged or older adults, are having the toughest time in today's economy. An analysis of government economic data suggests that this perception is correct. The recent indicators on the nation's labor market show a decline in the

Tough economic times altering young adults' daily lives, long-term plans. While negative trends in the labor market have been felt most acutely by the youngest workers, many adults in their late 20s and early 30s have also felt the impact of the weak economy. Among all 18- to 34-year-olds, fully half (49%) say they have taken a job they didn't want just to pay the bills, with 24% saying they have taken an unpaid job to gain work experience. And more than one-third (35%) say that, as a result of the poor economy, they have gone back to school. Their personal lives have also been affected: 31% have postponed either getting married or having a baby (22% say they have postponed having a baby and 20% have put off getting married). One-in-four (24%) say they have moved back in with their parents after living on their own.

http://pewresearch.org/pubs/2191/young-adults-workers-labor-market-pay-careers-advancement-recession

Margin of error

The general public survey is based on telephone interviews conducted Dec. 6-19, 2011, with a nationally representative sample of 2,048 adults ages 18 and older living in the continental United States, including an oversample of 346 adults ages 18 to 34. A total of 769 interviews were completed with respondents contacted by landline telephone and 1,279 with those contacted on their cellular phone. Data are weighted to produce a final sample that is representative of the general population of adults in the continental United States. Survey interviews were conducted under the direction of Princeton Survey Research Associates International, in English and Spanish. Margin of sampling error is plus or minus 2.9 percentage points for results based on the total sample and 4.4 percentage points for adults ages 18-34 at the 95% confidence level.

- 41% ± 2.9%: We are 95% confident that 38.1% to 43.9% of the public believe young adults, rather than middle-aged or older adults, are having the toughest time in today's economy.
- 49% ± 4.4%: We are 95% confident that 44.6% to 53.4% of 18-34 years olds have taken a job they didn't want just to pay the bills.

Parameter estimation

- We are often interested in *population parameters*.
- Since complete populations are difficult (or impossible) to collect data on, we use *sample statistics* as *point estimates* for the unknown population parameters of interest.
- Sample statistics vary from sample to sample.
- Quantifying how sample statistics vary provides a way to estimate the margin of error associated with our point estimate.
- But before we get to quantifying the variability among samples, let's try to understand how and why point estimates vary from sample to sample.

Suppose we randomly sample 1,000 adults from each state in the US. Would you expect the sample means of their heights to be the same, somewhat different, or very different?

Parameter estimation

- We are often interested in *population parameters*.
- Since complete populations are difficult (or impossible) to collect data on, we use *sample statistics* as *point estimates* for the unknown population parameters of interest.
- Sample statistics vary from sample to sample.
- Quantifying how sample statistics vary provides a way to estimate the margin of error associated with our point estimate.
- But before we get to quantifying the variability among samples, let's try to understand how and why point estimates vary from sample to sample.

Suppose we randomly sample 1,000 adults from each state in the US. Would you expect the sample means of their heights to be the same, somewhat different, or very different?

Not the same, but only somewhat different.

The following histogram shows the distribution of number of drinks it takes a group of college students to get drunk. We will assume that this is our population of interest. If we randomly select observations from this data set, which values are most likely to be selected, which are least likely?



Number of drinks to get drunk

Suppose that you don't have access to the population data. In order to estimate the average number of drinks it takes these college students to get drunk, you might sample from the population and use your sample mean as the best guess for the unknown population mean.

- Sample, with replacement, ten students from the population, and record the number of drinks it takes them to get drunk.
- Find the sample mean.
- Plot the distribution of the sample averages obtained by members of the class.

1	7	16	3	31	5	46	4	61	10	76	6	91	4	106	6	121	6	136	6
2	5	17	10	32	9	47	3	62	7	77	6	92	0.5	107	2	122	5	137	7
3	4	18	8	33	7	48	3	63	4	78	5	93	3	108	5	123	3	138	3
4	4	19	5	34	5	49	6	64	5	79	4	94	3	109	1	124	2	139	10
5	6	20	10	35	5	50	8	65	6	80	5	95	5	110	5	125	2	140	4
6	2	21	6	36	7	51	8	66	6	81	6	96	6	111	5	126	5	141	4
7	3	22	2	37	4	52	8	67	6	82	5	97	4	112	4	127	10	142	6
8	5	23	6	38	0	53	2	68	7	83	6	98	4	113	4	128	4	143	6
9	5	24	7	39	4	54	4	69	7	84	8	99	2	114	9	129	1	144	4
10	6	25	3	40	3	55	8	70	5	85	4	100	5	115	4	130	4	145	5
11	1	26	6	41	6	56	3	71	10	86	10	101	4	116	3	131	10	146	5
12	10	27	5	42	10	57	5	72	3	87	5	102	7	117	3	132	8		
13	4	28	8	43	3	58	5	73	5.5	88	10	103	6	118	4	133	10		
14	4	29	0	44	6	59	8	74	7	89	8	104	8	119	4	134	6		
15	6	30	8	45	10	60	4	75	10	90	5	105	3	120	8	135	6		

Example:

List of random numbers: 59, 121, 88, 46, 58, 72, 82, 81, 5, 10

1	7	16	3	31	5	46	4	61	10	76	6	91	4	106	6	121	6	136	6
2	5	17	10	32	9	47	3	62	7	77	6	92	0.5	107	2	122	5	137	7
3	4	18	8	33	7	48	3	63	4	78	5	93	3	108	5	123	3	138	3
4	4	19	5	34	5	49	6	64	5	79	4	94	3	109	1	124	2	139	10
5	6	20	10	35	5	50	8	65	6	80	5	95	5	110	5	125	2	140	4
6	2	21	6	36	7	51	8	66	6	81	6	96	6	111	5	126	5	141	4
7	3	22	2	37	4	52	8	67	6	82	5	97	4	112	4	127	10	142	6
8	5	23	6	38	0	53	2	68	7	83	6	98	4	113	4	128	4	143	6
9	5	24	7	39	4	54	4	69	7	84	8	99	2	114	9	129	1	144	4
10	6	25	3	40	3	55	8	70	5	85	4	100	5	115	4	130	4	145	5
11	1	26	6	41	6	56	3	71	10	86	10	101	4	116	3	131	10	146	5
12	10	27	5	42	10	57	5	72	3	87	5	102	7	117	3	132	8		
13	4	28	8	43	3	58	5	73	5.5	88	10	103	6	118	4	133	10		
14	4	29	0	44	6	59	8	74	7	89	8	104	8	119	4	134	6		
15	6	30	8	45	10	60	4	75	10	90	5	105	3	120	8	135	6		

Example:

List of random numbers: 59, 121, 88, 46, 58, 72, 82, 81, 5, 10

1	7	16	3	31	5	46	4	61	10	76	6	91	4	106	6	121	6	136	6
2	5	17	10	32	9	47	3	62	7	77	6	92	0.5	107	2	122	5	137	7
3	4	18	8	33	7	48	3	63	4	78	5	93	3	108	5	123	3	138	3
4	4	19	5	34	5	49	6	64	5	79	4	94	3	109	1	124	2	139	10
5	6	20	10	35	5	50	8	65	6	80	5	95	5	110	5	125	2	140	4
6	2	21	6	36	7	51	8	66	6	81	6	96	6	111	5	126	5	141	4
7	3	22	2	37	4	52	8	67	6	82	5	97	4	112	4	127	10	142	6
8	5	23	6	38	0	53	2	68	7	83	6	98	4	113	4	128	4	143	6
9	5	24	7	39	4	54	4	69	7	84	8	99	2	114	9	129	1	144	4
10	6	25	3	40	3	55	8	70	5	85	4	100	5	115	4	130	4	145	5
11	1	26	6	41	6	56	3	71	10	86	10	101	4	116	3	131	10	146	5
12	10	27	5	42	10	57	5	72	3	87	5	102	7	117	З	132	8		
13	4	28	8	43	3	58	5	73	5.5	88	10	103	6	118	4	133	10		
14	4	29	0	44	6	59	8	74	7	89	8	104	8	119	4	134	6		
15	6	30	8	45	10	60	4	75	10	90	5	105	3	120	8	135	6		

Sample mean: (8+6+10+4+5+3+5+6+6+6) / 10 = 5.9

Sampling distribution

What you just constructed is called a *sampling distribution*.

Sampling distribution

What you just constructed is called a *sampling distribution*.

What is the shape and center of this distribution? Based on this distribution, what do you think is the true population average?

Sampling distribution

What you just constructed is called a *sampling distribution*.

What is the shape and center of this distribution? Based on this distribution, what do you think is the true population average?

Approximately 5.39, the true population mean.

Central limit theorem

Central limit theorem

The distribution of the sample mean is well approximated by a normal model:

$$\bar{x} \sim N\left(mean = \mu, SE = \frac{\sigma}{\sqrt{n}}\right),$$

where SE is represents *standard error*, which is defined as the standard deviation of the sampling distribution. If σ is unknown, use *s*.

- It wasn't a coincidence that the sampling distribution we saw earlier was symmetric, and centered at the true population mean.
- We won't go through a detailed proof of why $SE = \frac{\sigma}{\sqrt{n}}$, but note that as *n* increases *SE* decreases.
 - As the sample size increases we would expect samples to yield more consistent sample means, hence the variability among the sample means would be lower.

CLT - conditions

Certain conditions must be met for the CLT to apply:

- 1. *Independence:* Sampled observations must be independent. This is difficult to verify, but is more likely if
 - random sampling/assignment is used, and
 - if sampling without replacement, n < 10% of the population.

CLT - conditions

Certain conditions must be met for the CLT to apply:

- 1. *Independence:* Sampled observations must be independent. This is difficult to verify, but is more likely if
 - random sampling/assignment is used, and
 - if sampling without replacement, n < 10% of the population.
- 2. *Sample size/skew:* Either the population distribution is normal, or if the population distribution is skewed, the sample size is large.
 - the more skewed the population distribution, the larger sample size we need for the CLT to apply
 - for moderately skewed distributions n > 30 is a widely used rule of thumb

This is also difficult to verify for the population, but we can check it using the sample data, and assume that the sample mirrors the population.

1) Variability in estimates

2 Confidence intervals

- Why do we report confidence intervals?
- Constructing a confidence interval
- A more accurate interval
- Capturing the population parameter
- Changing the confidence level

3 Hypothesis testing

- Examining the Central Limit Theorem
- 5 Inference for other estimators
- 6 Sample size and power



Confidence intervals

- A plausible range of values for the population parameter is called a *confidence interval*.
- Using only a sample statistic to estimate a parameter is like fishing in a murky lake with a spear, and using a confidence interval is like fishing with a net.



We can throw a spear where we saw a fish but we will probably miss. If we toss a net in that area, we have a good chance of catching the fish.



 If we report a point estimate, we probably won't hit the exact population parameter. If we report a range of plausible values we have a good shot at capturing the parameter.

Photos by Mark Fischer (http://www.flickr.com/photos/fischerfotos/7439791462) and Chris Penny

(http://www.flickr.com/photos/clearlydived/7029109617) on Flickr.

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

The approximate 95% confidence interval is defined as

point estimate $\pm 2 \times SE$

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

The approximate 95% confidence interval is defined as

point estimate $\pm 2 \times SE$

$$SE = \frac{s}{\sqrt{n}} = \frac{1.74}{\sqrt{50}} \approx 0.25$$

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

The approximate 95% confidence interval is defined as

point estimate $\pm 2 \times SE$

$$SE = \frac{s}{\sqrt{n}} = \frac{1.74}{\sqrt{50}} \approx 0.25$$

 $\bar{x} \pm 2 \times SE = 3.2 \pm 2 \times 0.25$

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

The approximate 95% confidence interval is defined as

point estimate $\pm 2 \times SE$

$$SE = \frac{s}{\sqrt{n}} = \frac{1.74}{\sqrt{50}} \approx 0.25$$

$$\bar{x} \pm 2 \times SE = 3.2 \pm 2 \times 0.25$$

= (3.2 - 0.5, 3.2 + 0.5)

A random sample of 50 college students were asked how many exclusive relationships they have been in so far. This sample yielded a mean of 3.2 and a standard deviation of 1.74. Estimate the true average number of exclusive relationships using this sample.

$$\bar{x} = 3.2$$
 $s = 1.74$

The approximate 95% confidence interval is defined as

point estimate $\pm 2 \times SE$

$$SE = \frac{s}{\sqrt{n}} = \frac{1.74}{\sqrt{50}} \approx 0.25$$

$$\bar{x} \pm 2 \times SE = 3.2 \pm 2 \times 0.25$$

= (3.2 - 0.5, 3.2 + 0.5)
= (2.7, 3.7)

Which of the following is the correct interpretation of this confidence interval?

We are 95% confident that

- (a) the average number of exclusive relationships college students in this sample have been in is between 2.7 and 3.7.
- (b) college students on average have been in between 2.7 and 3.7 exclusive relationships.
- (c) a randomly chosen college student has been in 2.7 to 3.7 exclusive relationships.
- (d) 95% of college students have been in 2.7 to 3.7 exclusive relationships.

Which of the following is the correct interpretation of this confidence interval?

We are 95% confident that

- (a) the average number of exclusive relationships college students in this sample have been in is between 2.7 and 3.7.
- (b) college students on average have been in between 2.7 and 3.7 exclusive relationships.
- (c) a randomly chosen college student has been in 2.7 to 3.7 exclusive relationships.
- (d) 95% of college students have been in 2.7 to 3.7 exclusive relationships.

A more accurate interval

Confidence interval, a general formula

point estimate $\pm z^* \times SE$

A more accurate interval

Confidence interval, a general formula

point estimate $\pm z^* \times SE$

Conditions when the point estimate = \bar{x} :

- 1. Independence: Observations in the sample must be independent
 - random sample/assignment
 - if sampling without replacement, n < 10% of population
- 2. Sample size / skew: $n \ge 30$ and population distribution should not be extremely skewed

A more accurate interval

Confidence interval, a general formula

point estimate $\pm z^{\star} \times SE$

Conditions when the point estimate = \bar{x} :

1. Independence: Observations in the sample must be independent

- random sample/assignment
- if sampling without replacement, n < 10% of population
- 2. Sample size / skew: $n \ge 30$ and population distribution should not be extremely skewed

Note: We will discuss working with samples where n < 30 in the next chapter.

OpenIntro Statistics, 2nd Edition

What does 95% confident mean?

- Suppose we took many samples and built a confidence interval from each sample using the equation *point estimate* $\pm 2 \times SE$.
- Then about 95% of those intervals would contain the true population mean (μ).
- The figure shows this process with 25 samples, where 24 of the resulting confidence intervals contain the true average number of exclusive relationships, and one does not.



If we want to be more certain that we capture the population parameter, i.e. increase our confidence level, should we use a wider interval or a smaller interval?

If we want to be more certain that we capture the population parameter, i.e. increase our confidence level, should we use a wider interval or a smaller interval?

A wider interval.

If we want to be more certain that we capture the population parameter, i.e. increase our confidence level, should we use a wider interval or a smaller interval?

A wider interval.

Can you see any drawbacks to using a wider interval?



If we want to be more certain that we capture the population parameter, i.e. increase our confidence level, should we use a wider interval or a smaller interval?

A wider interval.

Can you see any drawbacks to using a wider interval?



If the interval is too wide it may not be very informative.

OpenIntro Statistics, 2nd Edition
Image source: http://web.as.uky.edu/statistics/users/earo227/misc/garfield_weather.gif

Changing the confidence level

point estimate $\pm z^* \times SE$

- In a confidence interval, *z** × *SE* is called the *margin of error*, and for a given sample, the margin of error changes as the confidence level changes.
- In order to change the confidence level we need to adjust z* in the above formula.
- Commonly used confidence levels in practice are 90%, 95%, 98%, and 99%.
- For a 95% confidence interval, $z^* = 1.96$.
- However, using the standard normal (z) distribution, it is possible to find the appropriate z^* for any confidence level.

Which of the below Z scores is the appropriate z^* when calculating a 98% confidence interval?

- (a) Z = 2.05 (d) Z = -2.33
- (b) Z = 1.96 (e) Z = -1.65
- (c) Z = 2.33

Which of the below Z scores is the appropriate z^* when calculating a 98% confidence interval?

- (a) Z = 2.05 (d) Z = -2.33
- (b) Z = 1.96 (e) Z = -1.65
- (c) Z = 2.33



Variability in estimates

2 Confidence intervals

3 Hypothesis testing

- Hypothesis testing framework
- Testing hypotheses using confidence intervals
- Conditions for inference
- Formal testing using p-values
- Two-sided hypothesis testing with p-values
- Decision errors
- Choosing a significance level
- Recap



- 5 Inference for other estimators
- Sample size and power



Remember when...

Gender di	scriminati	on experime	ent:		
		Promotion			
		Promoted	Not Promoted	Total	
Gender	Male	21	3	24	
	Female	14	10	24	
	Total	35	13	48	

Remember when...

Condor discrimination ovporimont:

Gender discrimination experiment.						
		Promotion				
		Promoted	Not Promoted	Total		
Gender	Male	21	3	24		
	Female	14	10	24		
	Total	35	13	48		
	$\hat{p}_{males} = 21/24 \approx 0.88$					
	$\hat{p}_{females} = 14/24 \approx 0.58$					

Remember when...

Gender discrimination experiment:							
		Pro					
		Promoted	Not Promoted	Total			
Gender	Male	21	3	24			
	Female	14	10	24			
	Total	35	13	48			
	$\hat{p}_{males} = 21/24 \approx 0.88$						
$\hat{p}_{females} = 14/24 \approx 0.58$							

Possible explanations:

- Promotion and gender are *independent*, no gender discrimination, observed difference in proportions is simply due to chance. → *null* - (nothing is going on)
- Promotion and gender are *dependent*, there is gender discrimination, observed difference in proportions is not due to chance. → *alternative* - (something is going on)

Result



Difference in promotion rates

Result



Difference in promotion rates

Since it was quite unlikely to obtain results like the actual data or something more extreme in the simulations (male promotions being 30% or more higher than female promotions), we decided to reject the null hypothesis in favor of the alternative.

OpenIntro Statistics, 2nd Edition

• We start with a *null hypothesis* (*H*₀) that represents the status quo.

- We start with a *null hypothesis* (*H*₀) that represents the status quo.
- We also have an *alternative hypothesis* (*H_A*) that represents our research question, i.e. what we're testing for.

- We start with a *null hypothesis* (*H*₀) that represents the status quo.
- We also have an *alternative hypothesis* (*H_A*) that represents our research question, i.e. what we're testing for.
- We conduct a hypothesis test under the assumption that the null hypothesis is true, either via simulation or traditional methods based on the central limit theorem (coming up next...).

- We start with a *null hypothesis* (*H*₀) that represents the status quo.
- We also have an *alternative hypothesis* (*H_A*) that represents our research question, i.e. what we're testing for.
- We conduct a hypothesis test under the assumption that the null hypothesis is true, either via simulation or traditional methods based on the central limit theorem (coming up next...).
- If the test results suggest that the data do not provide convincing evidence for the alternative hypothesis, we stick with the null hypothesis. If they do, then we reject the null hypothesis in favor of the alternative.

- We start with a *null hypothesis* (*H*₀) that represents the status quo.
- We also have an *alternative hypothesis* (*H_A*) that represents our research question, i.e. what we're testing for.
- We conduct a hypothesis test under the assumption that the null hypothesis is true, either via simulation or traditional methods based on the central limit theorem (coming up next...).
- If the test results suggest that the data do not provide convincing evidence for the alternative hypothesis, we stick with the null hypothesis. If they do, then we reject the null hypothesis in favor of the alternative.

We'll formally introduce the hypothesis testing framework using an example on testing a claim about a population mean.

Earlier we calculated a 95% confidence interval for the average number of exclusive relationships college students have been in to be (2.7, 3.7). Based on this confidence interval, do these data support the hypothesis that college students on average have been in more than 3 exclusive relationships.

Earlier we calculated a 95% confidence interval for the average number of exclusive relationships college students have been in to be (2.7, 3.7). Based on this confidence interval, do these data support the hypothesis that college students on average have been in more than 3 exclusive relationships.

- The associated hypotheses are:
 - *H*₀: μ = 3: College students have been in 3 exclusive relationships, on average
 - *H_A*: $\mu > 3$: College students have been in more than 3 exclusive relationships, on average

Earlier we calculated a 95% confidence interval for the average number of exclusive relationships college students have been in to be (2.7, 3.7). Based on this confidence interval, do these data support the hypothesis that college students on average have been in more than 3 exclusive relationships.

- The associated hypotheses are:
 - *H*₀: μ = 3: College students have been in 3 exclusive relationships, on average
 - *H_A*: $\mu > 3$: College students have been in more than 3 exclusive relationships, on average
- Since the null value is included in the interval, we do not reject the null hypothesis in favor of the alternative.

Earlier we calculated a 95% confidence interval for the average number of exclusive relationships college students have been in to be (2.7, 3.7). Based on this confidence interval, do these data support the hypothesis that college students on average have been in more than 3 exclusive relationships.

- The associated hypotheses are:
 - *H*₀: μ = 3: College students have been in 3 exclusive relationships, on average
 - *H_A*: $\mu > 3$: College students have been in more than 3 exclusive relationships, on average
- Since the null value is included in the interval, we do not reject the null hypothesis in favor of the alternative.
- This is a quick-and-dirty approach for hypothesis testing. However it doesn't tell us the likelihood of certain outcomes under the null hypothesis, i.e. the p-value, based on which we can make a decision on the hypotheses.

OpenIntro Statistics, 2nd Edition

Chp 4: Foundations for inference

Number of college applications

A similar survey asked how many colleges students applied to, and 206 students responded to this question. This sample yielded an average of 9.7 college applications with a standard deviation of 7. College Board website states that counselors recommend students apply to roughly 8 colleges. Do these data provide convincing evidence that the average number of colleges all Duke students apply to is <u>higher</u> than recommended?

http://www.collegeboard.com/student/apply/the-application/151680.html

• The *parameter of interest* is the average number of schools applied to by all Duke students.

- The *parameter of interest* is the average number of schools applied to by <u>all</u> Duke students.
- There may be two explanations why our sample mean is higher than the recommended 8 schools.
 - The true population mean is different.
 - The true population mean is 8, and the difference between the true population mean and the sample mean is simply due to natural sampling variability.

- The *parameter of interest* is the average number of schools applied to by <u>all</u> Duke students.
- There may be two explanations why our sample mean is higher than the recommended 8 schools.
 - The true population mean is different.
 - The true population mean is 8, and the difference between the true population mean and the sample mean is simply due to natural sampling variability.
- We start with the assumption the average number of colleges Duke students apply to is 8 (as recommended)

 $H_0: \mu = 8$

- The *parameter of interest* is the average number of schools applied to by <u>all</u> Duke students.
- There may be two explanations why our sample mean is higher than the recommended 8 schools.
 - The true population mean is different.
 - The true population mean is 8, and the difference between the true population mean and the sample mean is simply due to natural sampling variability.
- We start with the assumption the average number of colleges Duke students apply to is 8 (as recommended)

$$H_0: \mu = 8$$

• We test the claim that the average number of colleges Duke students apply to is greater than 8

$$H_A: \mu > 8$$

Number of college applications - conditions

Which of the following is <u>not</u> a condition that needs to be met to proceed with this hypothesis test?

- (a) Students in the sample should be independent of each other with respect to how many colleges they applied to.
- (b) Sampling should have been done randomly.
- (c) The sample size should be less than 10% of the population of all Duke students.
- (d) There should be at least 10 successes and 10 failures in the sample.
- (e) The distribution of the number of colleges students apply to should not be extremely skewed.

Number of college applications - conditions

Which of the following is <u>not</u> a condition that needs to be met to proceed with this hypothesis test?

- (a) Students in the sample should be independent of each other with respect to how many colleges they applied to.
- (b) Sampling should have been done randomly.
- (c) The sample size should be less than 10% of the population of all Duke students.
- (d) There should be at least 10 successes and 10 failures in the sample.
- (e) The distribution of the number of colleges students apply to should not be extremely skewed.





$$\bar{x} \sim N\left(\mu = 8, SE = \frac{7}{\sqrt{206}} = 0.5\right)$$



$$\bar{x} \sim N\left(\mu = 8, SE = \frac{7}{\sqrt{206}} = 0.5\right)$$

$$Z = \frac{9.7 - 8}{0.5} = 3.4$$

In order to evaluate if the observed sample mean is unusual for the hypothesized sampling distribution, we determine how many standard errors away from the null it is, which is also called the *test statistic*.



The sample mean is 3.4 standard errors away from the hypothesized value. Is this considered unusually high? That is, is the result *statistically significant*?

In order to evaluate if the observed sample mean is unusual for the hypothesized sampling distribution, we determine how many standard errors away from the null it is, which is also called the *test statistic*.



The sample mean is 3.4 standard errors away from the hypothesized value. Is this considered unusually high? That is, is the result *statistically significant*?

Yes, and we can quantify how unusual it is using a p-value.

p-values

 We then use this test statistic to calculate the *p-value*, the probability of observing data at least as favorable to the alternative hypothesis as our current data set, if the null hypothesis were true.

p-values

- We then use this test statistic to calculate the *p-value*, the probability of observing data at least as favorable to the alternative hypothesis as our current data set, if the null hypothesis were true.
- If the p-value is *low* (lower than the significance level, α, which is usually 5%) we say that it would be very unlikely to observe the data if the null hypothesis were true, and hence *reject* H₀.

p-values

- We then use this test statistic to calculate the *p-value*, the probability of observing data at least as favorable to the alternative hypothesis as our current data set, if the null hypothesis were true.
- If the p-value is *low* (lower than the significance level, α, which is usually 5%) we say that it would be very unlikely to observe the data if the null hypothesis were true, and hence *reject* H₀.
- If the p-value is *high* (higher than α) we say that it is likely to observe the data even if the null hypothesis were true, and hence *do not reject H*₀.

Number of college applications - p-value

p-value: probability of observing data at least as favorable to H_A as our current data set (a sample mean greater than 9.7), if in fact H_0 were true (the true population mean was 8).



Number of college applications - p-value

p-value: probability of observing data at least as favorable to H_A as our current data set (a sample mean greater than 9.7), if in fact H_0 were true (the true population mean was 8).



 $P(\bar{x} > 9.7 \mid \mu = 8) = P(Z > 3.4) = 0.0003$
• p-value = 0.0003

- p-value = 0.0003
 - If the true average of the number of colleges Duke students applied to is 8, there is only 0.03% chance of observing a random sample of 206 Duke students who on average apply to 9.7 or more schools.

• p-value = 0.0003

- If the true average of the number of colleges Duke students applied to is 8, there is only 0.03% chance of observing a random sample of 206 Duke students who on average apply to 9.7 or more schools.
- This is a pretty low probability for us to think that a sample mean of 9.7 or more schools is likely to happen simply by chance.

- p-value = 0.0003
 - If the true average of the number of colleges Duke students applied to is 8, there is only 0.03% chance of observing a random sample of 206 Duke students who on average apply to 9.7 or more schools.
 - This is a pretty low probability for us to think that a sample mean of 9.7 or more schools is likely to happen simply by chance.
- Since p-value is *low* (lower than 5%) we *reject* H_0 .

- p-value = 0.0003
 - If the true average of the number of colleges Duke students applied to is 8, there is only 0.03% chance of observing a random sample of 206 Duke students who on average apply to 9.7 or more schools.
 - This is a pretty low probability for us to think that a sample mean of 9.7 or more schools is likely to happen simply by chance.
- Since p-value is *low* (lower than 5%) we *reject* H_0 .
- The data provide convincing evidence that Duke students apply to more than 8 schools on average.

- p-value = 0.0003
 - If the true average of the number of colleges Duke students applied to is 8, there is only 0.03% chance of observing a random sample of 206 Duke students who on average apply to 9.7 or more schools.
 - This is a pretty low probability for us to think that a sample mean of 9.7 or more schools is likely to happen simply by chance.
- Since p-value is *low* (lower than 5%) we *reject* H_0 .
- The data provide convincing evidence that Duke students apply to more than 8 schools on average.
- The difference between the null value of 8 schools and observed sample mean of 9.7 schools is *not due to chance* or sampling variability.

A poll by the National Sleep Foundation found that college students average about 7 hours of sleep per night. A sample of 169 college students taking an introductory statistics class yielded an average of 6.88 hours, with a standard deviation of 0.94 hours. Assuming that this is a random sample representative of all college students *(bit of a leap of faith?)*, a hypothesis test was conducted to evaluate if college students on average sleep less than 7 hours per night. The p-value for this hypothesis test is 0.0485. Which of the following is correct?

- (a) Fail to reject H_0 , the data provide convincing evidence that college students sleep less than 7 hours on average.
- (b) Reject H_0 , the data provide convincing evidence that college students sleep less than 7 hours on average.
- (c) Reject H_0 , the data prove that college students sleep more than 7 hours on average.
- (d) Fail to reject H_0 , the data do not provide convincing evidence that college students sleep less than 7 hours on average.
- (e) Reject H_0 , the data provide convincing evidence that college students in this sample sleep less than 7 hours on average.

A poll by the National Sleep Foundation found that college students average about 7 hours of sleep per night. A sample of 169 college students taking an introductory statistics class yielded an average of 6.88 hours, with a standard deviation of 0.94 hours. Assuming that this is a random sample representative of all college students *(bit of a leap of faith?)*, a hypothesis test was conducted to evaluate if college students on average sleep less than 7 hours per night. The p-value for this hypothesis test is 0.0485. Which of the following is correct?

- (a) Fail to reject H_0 , the data provide convincing evidence that college students sleep less than 7 hours on average.
- (b) Reject H₀, the data provide convincing evidence that college students sleep less than 7 hours on average.
- (c) Reject H_0 , the data prove that college students sleep more than 7 hours on average.
- (d) Fail to reject H_0 , the data do not provide convincing evidence that college students sleep less than 7 hours on average.
- (e) Reject H_0 , the data provide convincing evidence that college students in this sample sleep less than 7 hours on average.

OpenIntro Statistics, 2nd Edition

Two-sided hypothesis testing with p-values

 If the research question was "Do the data provide convincing evidence that the average amount of sleep college students get per night is *different* than the national average?", the alternative hypothesis would be different.

 $H_0: \mu = 7$ $H_A: \mu \neq 7$

Two-sided hypothesis testing with p-values

 If the research question was "Do the data provide convincing evidence that the average amount of sleep college students get per night is *different* than the national average?", the alternative hypothesis would be different.

 $H_0: \mu = 7$ $H_A: \mu \neq 7$

Hence the p-value would change as well:



Decision errors

- Hypothesis tests are not flawless.
- In the court system innocent people are sometimes wrongly convicted and the guilty sometimes walk free.
- Similarly, we can make a wrong decision in statistical hypothesis tests as well.
- The difference is that we have the tools necessary to quantify how often we make errors in statistics.

		Decision	
		fail to reject H_0	reject H_0
Turala	H ₀ true		
Iruth	H_A true		

		Decision	
		fail to reject H_0	reject H_0
T	H ₀ true	\checkmark	
Iruth	H_A true		



There are two competing hypotheses: the null and the alternative. In a hypothesis test, we make a decision about which might be true, but our choice might be incorrect.

		Decision	
		fail to reject H_0	reject H_0
Truth	H ₀ true	\checkmark	Type 1 Error
	H_A true		\checkmark

• A *Type 1 Error* is rejecting the null hypothesis when H_0 is true.

		Decision	
		fail to reject H_0	reject H_0
T	H ₀ true	\checkmark	Type 1 Error
Iruth	H_A true	Type 2 Error	\checkmark

- A Type 1 Error is rejecting the null hypothesis when H_0 is true.
- A *Type 2 Error* is failing to reject the null hypothesis when *H*_A is true.

		Decision	
		fail to reject H_0	reject H_0
Turala	H ₀ true	\checkmark	Type 1 Error
Iruth	H_A true	Type 2 Error	\checkmark

- A Type 1 Error is rejecting the null hypothesis when H_0 is true.
- A *Type 2 Error* is failing to reject the null hypothesis when *H*_A is true.
- We (almost) never know if *H*₀ or *H*_A is true, but we need to consider all possibilities.

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

- H_0 : Defendant is innocent
- H_A : Defendant is guilty

Which type of error is being committed in the following cirumstances?

- Declaring the defendant innocent when they are actually guilty
- Declaring the defendant guilty when they are actually innocent

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

- H_0 : Defendant is innocent
- H_A : Defendant is guilty

Which type of error is being committed in the following cirumstances?

• Declaring the defendant innocent when they are actually guilty

Type 2 error

Declaring the defendant guilty when they are actually innocent

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

- H_0 : Defendant is innocent
- H_A : Defendant is guilty

Which type of error is being committed in the following cirumstances?

• Declaring the defendant innocent when they are actually guilty

Type 2 error

• Declaring the defendant guilty when they are actually innocent Type 1 error

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

- H_0 : Defendant is innocent
- H_A : Defendant is guilty

Which type of error is being committed in the following cirumstances?

• Declaring the defendant innocent when they are actually guilty

Type 2 error

• Declaring the defendant guilty when they are actually innocent *Type 1 error*

Which error do you think is the worse error to make?

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

- H_0 : Defendant is innocent
- H_A : Defendant is guilty

Which type of error is being committed in the following cirumstances?

- Declaring the defendant innocent when they are actually guilty Type 2 error
- Declaring the defendant guilty when they are actually innocent *Type 1 error*

Which error do you think is the worse error to make?

"better that ten guilty persons escape than that one innocent suffer"

- William Blackstone

• As a general rule we reject H_0 when the p-value is less than 0.05, i.e. we use a *significance level* of 0.05, $\alpha = 0.05$.

- As a general rule we reject H_0 when the p-value is less than 0.05, i.e. we use a *significance level* of 0.05, $\alpha = 0.05$.
- This means that, for those cases where *H*₀ is actually true, we do not want to incorrectly reject it more than 5% of those times.

- As a general rule we reject H_0 when the p-value is less than 0.05, i.e. we use a *significance level* of 0.05, $\alpha = 0.05$.
- This means that, for those cases where *H*₀ is actually true, we do not want to incorrectly reject it more than 5% of those times.
- In other words, when using a 5% significance level there is about 5% chance of making a Type 1 error if the null hypothesis is true.

 $P(\text{Type 1 error}) = \alpha$

- As a general rule we reject H_0 when the p-value is less than 0.05, i.e. we use a *significance level* of 0.05, $\alpha = 0.05$.
- This means that, for those cases where *H*₀ is actually true, we do not want to incorrectly reject it more than 5% of those times.
- In other words, when using a 5% significance level there is about 5% chance of making a Type 1 error if the null hypothesis is true.

 $P(\text{Type 1 error}) = \alpha$

 This is why we prefer small values of α – increasing α increases the Type 1 error rate.

Choosing a significance level

- Choosing a significance level for a test is important in many contexts, and the traditional level is 0.05. However, it is often helpful to adjust the significance level based on the application.
- We may select a level that is smaller or larger than 0.05 depending on the consequences of any conclusions reached from the test.
- If making a Type 1 Error is dangerous or especially costly, we should choose a small significance level (e.g. 0.01). Under this scenario we want to be very cautious about rejecting the null hypothesis, so we demand very strong evidence favoring H_A before we would reject H_0 .
- If a Type 2 Error is relatively more dangerous or much more costly than a Type 1 Error, then we should choose a higher significance level (e.g. 0.10). Here we want to be cautious about failing to reject H₀ when the null is actually false.

the next two slides are provided as a brief summary of hypothesis testing ...

Recap

Recap: Hypothesis testing framework

- 1. Set the hypotheses.
- 2. Check assumptions and conditions.
- 3. Calculate a *test statistic* and a p-value.
- 4. Make a decision, and interpret it in context of the research question.

Recap: Hypothesis testing for a population mean

- 1. Set the hypotheses
 - $H_0: \mu = null value$
 - $H_A: \mu < \text{or} > \text{or} \neq null value}$
- Calculate the point estimate
- Check assumptions and conditions
 - Independence: random sample/assignment, 10% condition when sampling without replacement
 - Normality: nearly normal population or $n \ge 30$, no extreme skew or use the t distribution
- 4. Calculate a *test statistic* and a p-value (draw a picture!)

$$Z = \frac{\bar{x} - \mu}{SE}$$
, where $SE = \frac{s}{\sqrt{n}}$

- Make a decision, and interpret it in context
 - If p-value $< \alpha$, reject H_0 , data provide evidence for H_A
 - If p-value > α , do not reject H_0 , data do not provide evidence for H_A

- Variability in estimates
- 2 Confidence intervals
- 3 Hypothesis testing
- Examining the Central Limit Theorem
- Inference for other estimators
- 6 Sample size and power
- Statistical vs. practical significance

Average number of basketball games attended

Next let's look at the population data for the number of basketball games attended:



Average number of basketball games attended (cont.)

Sampling distribution, n = 10:



What does each observation in this distribution represent?

Is the variability of the sampling distribution smaller or larger than the variability of the population distribution? Why?

Average number of basketball games attended (cont.)

Sampling distribution, n = 10:



What does each observation in this distribution represent?

Sample mean (\bar{x}) of samples of size n = 10.

Is the variability of the sampling distribution smaller or larger than the variability of the population distribution? Why?

Average number of basketball games attended (cont.)

Sampling distribution, n = 10:



sample means from samples of n = 10

What does each observation in this distribution represent?

Sample mean (\bar{x}) of samples of size n = 10.

Is the variability of the sampling distribution smaller or larger than the variability of the population distribution? Why?

Smaller, sample means will vary less than individual observations.




How did the shape, center, and spread of the sampling distribution change going from n = 10 to n = 30?





How did the shape, center, and spread of the sampling distribution change going from n = 10 to n = 30?

Shape is more symmetric, center is about the same, spread is smaller.

Sampling distribution, n = 70:



The mean of the sampling distribution is 5.75, and the standard deviation of the sampling distribution (also called the *standard error*) is 0.75. Which of the following is the most reasonable guess for the 95% confidence interval for the true average number of basketball games attended by students?

- (a) 5.75 ± 0.75
- (b) $5.75 \pm 2 \times 0.75$
- (c) $5.75 \pm 3 \times 0.75$
- (d) cannot tell from the information given

The mean of the sampling distribution is 5.75, and the standard deviation of the sampling distribution (also called the *standard error*) is 0.75. Which of the following is the most reasonable guess for the 95% confidence interval for the true average number of basketball games attended by students?

- (a) 5.75 ± 0.75
- (b) $5.75 \pm 2 \times 0.75 \rightarrow (4.25, 7.25)$
- (c) $5.75 \pm 3 \times 0.75$
- (d) cannot tell from the information given

Four plots: Determine which plot (A, B, or C) is which.

- (1) At top: distribution for a population ($\mu = 10, \sigma = 7$),
- (2) a single random sample of 100 observations from this population,
- (3) a distribution of 100 sample means from random samples with size 7, and

(4) a distribution of 100 sample means from random samples with size 49.



Four plots: Determine which plot (A, B, or C) is which.

- (1) At top: distribution for a population ($\mu = 10, \sigma = 7$),
- (2) a single random sample of 100 observations from this population,
- (3) a distribution of 100 sample means from random samples with size 7, and

(4) a distribution of 100 sample means from random samples with size 49.



- Variability in estimates
- 2 Confidence intervals
- 3 Hypothesis testing
- Examining the Central Limit Theorem

5 Inference for other estimators

- Confidence intervals for nearly normal point estimates
- Hypothesis testing for nearly normal point estimates
- Non-normal point estimates
- When to retreat

6 Sample size and power

Statistical vs. practical significance

Inference for other estimators

- The sample mean is not the only point estimate for which the sampling distribution is nearly normal. For example, the sampling distribution of sample proportions is also nearly normal when *n* is sufficiently large.
- An important assumption about point estimates is that they are *unbiased*, i.e. the sampling distribution of the estimate is centered at the true population parameter it estimates.
 - That is, an unbiased estimate does not naturally over or underestimate the parameter. Rather, it tends to provide a "good" estimate.
 - The sample mean is an example of an unbiased point estimate, as are each of the examples we introduce in this section.
- Some point estimates follow distributions other than the normal distribution, and some scenarios require statistical techniques that we havenÕt covered yet – we will discuss these at the end of this section.

Confidence intervals for nearly normal point estimates

A confidence interval based on an unbiased and nearly normal point estimate is

point estimate $\pm z^*SE$

where z^* is selected to correspond to the confidence level, and SE represents the standard error.

Remember that the value z^*SE is called the *margin of error*.

One of the earliest examples of behavioral asymmetry is a preference in humans for turning the head to the right, rather than to the left, during the final weeks of gestation and for the first 6 months after birth. This is thought to influence subsequent development of perceptual and motor preferences. A study of 124 couples found that 64.5% turned their heads to the right when kissing. The standard error associated with this estimate is roughly 4%. Which of the below is false?

- (a) The 95% confidence interval for the percentage of kissers who turn their heads to the right is roughly $64.5\% \pm 4\%$.
- (b) A higher sample size would yield a lower standard error.
- (c) The margin of error for a 95% confidence interval for the percentage of kissers who turn their heads to the right is roughly 8%.
- (d) The 99.7% confidence interval for the percentage of kissers who turn their heads to the right is roughly $64.5\% \pm 12\%$.

Güntürkün, O. (2003) Adult persistence of head-turning asymmetry. Nature. Vol 421.

One of the earliest examples of behavioral asymmetry is a preference in humans for turning the head to the right, rather than to the left, during the final weeks of gestation and for the first 6 months after birth. This is thought to influence subsequent development of perceptual and motor preferences. A study of 124 couples found that 64.5% turned their heads to the right when kissing. The standard error associated with this estimate is roughly 4%. Which of the below is false?

- (a) The 95% confidence interval for the percentage of kissers who turn their heads to the right is roughly $64.5\% \pm 4\%$.
- (b) A higher sample size would yield a lower standard error.
- (c) The margin of error for a 95% confidence interval for the percentage of kissers who turn their heads to the right is roughly 8%.
- (d) The 99.7% confidence interval for the percentage of kissers who turn their heads to the right is roughly $64.5\% \pm 12\%$.

Güntürkün, O. (2003) Adult persistence of head-turning asymmetry. Nature. Vol 421.

Hypothesis testing for nearly normal point estimates

The third National Health and Nutrition Examination Survey collected body fat percentage (BF%) and gender data from 13,601 subjects ages 20 to 80. The average BF% for the 6,580 men in the sample was 23.9, and this value was 35.0 for the 7,021 women. The standard error for the difference between the average men and women BF%s was 0.114. Do these data provide convincing evidence that men and women have different average BF%s. You may assume that the distribution of the point estimate is nearly normal.

Hypothesis testing for nearly normal point estimates

The third National Health and Nutrition Examination Survey collected body fat percentage (BF%) and gender data from 13,601 subjects ages 20 to 80. The average BF% for the 6,580 men in the sample was 23.9, and this value was 35.0 for the 7,021 women. The standard error for the difference between the average men and women BF%s was 0.114. Do these data provide convincing evidence that men and women have different average BF%s. You may assume that the distribution of the point estimate is nearly normal.

- 1. Set hypotheses
- 2. Calculate point estimate
- 3. Check conditions
- 4. Draw sampling distribution, shade p-value
- 5. Calculate test statistics and p-value, make a decision

1. The null hypothesis is that men and women have equal average BF%, and the alternative is that these values are different.

 $H_0: \mu_{men} = \mu_{women} \qquad H_A: \mu_{men} \neq \mu_{women}$

1. The null hypothesis is that men and women have equal average BF%, and the alternative is that these values are different.

 $H_0: \mu_{men} = \mu_{women}$ $H_A: \mu_{men} \neq \mu_{women}$

2. The parameter of interest is the average difference in the population means of BF%s for men and women, and the point estimate for this parameter is the difference between the two sample means:

$$\bar{x}_{men} - \bar{x}_{women} = 23.9 - 35.0 = -11.1$$

1. The null hypothesis is that men and women have equal average BF%, and the alternative is that these values are different.

 $H_0: \mu_{men} = \mu_{women}$ $H_A: \mu_{men} \neq \mu_{women}$

2. The parameter of interest is the average difference in the population means of BF%s for men and women, and the point estimate for this parameter is the difference between the two sample means:

$$\bar{x}_{men} - \bar{x}_{women} = 23.9 - 35.0 = -11.1$$

3. We are assuming that the distribution of the point estimate is nearly normal (we will discuss details for checking this condition in the next chapter, however given the large sample sizes, the normality assumption doesn't seem unwarranted). 4. The sampling distribution will be centered at the null value $(\mu_{men} - \mu_{women} = 0)$, and the p-value is the area beyond the observed difference in sample means in both tails (lower than -11.1 and higher than 11.1).



 The test statistic is computed as the difference between the point estimate and the null value (-11.1 - 0 = -11.1), scaled by the standard error.

$$Z = \frac{11.1 - 0}{0.114} = 97.36$$

The Z score is huge! And hence the p-value will be tiny, allowing us to reject H_0 in favor of H_A .

5. The test statistic is computed as the difference between the point estimate and the null value (-11.1 - 0 = -11.1), scaled by the standard error.

$$Z = \frac{11.1 - 0}{0.114} = 97.36$$

The Z score is huge! And hence the p-value will be tiny, allowing us to reject H_0 in favor of H_A .

These data provide convincing evidence that the average BF% of men and women are different.

Non-normal point estimates

- We may apply the ideas of confidence intervals and hypothesis testing to cases where the point estimate or test statistic is not necessarily normal. There are many reasons why such a situation may arise:
 - the sample size is too small for the normal approximation to be valid;
 - the standard error estimate may be poor; or
 - the point estimate tends towards some distribution that is not the normal distribution.
- For each case where the normal approximation is not valid, our first task is always to understand and characterize the sampling distribution of the point estimate or test statistic. Next, we can apply the general frameworks for confidence intervals and hypothesis testing to these alternative distributions.

When to retreat

- Statistical tools rely on the following two main conditions:
 - Independence A random sample from less than 10% of the population ensures independence of observations. In experiments, this is ensured by random assignment. If independence fails, then advanced techniques must be used, and in some such cases, inference may not be possible.
 - Sample size and skew For example, if the sample size is too small, the skew too strong, or extreme outliers are present, then the normal model for the sample mean will fail.
- Whenever conditions are not satisfied for a statistical technique:
 - 1. Learn new methods that are appropriate for the data.
 - 2. Consult a statistician.
 - Ignore the failure of conditions. This last option effectively invalidates any analysis and may discredit novel and interesting findings.

- Variability in estimates
- 2 Confidence intervals
- 3 Hypothesis testing
- Examining the Central Limit Theorem
- 5 Inference for other estimators
- Sample size and power
 - Finding a sample size for a certain margin of error
 - Power and the Type 2 Error rate
 - Statistical vs. practical significance

We know that the critical value associated with the 96% confidence level: $z^* = 2.05$.

We know that the critical value associated with the 96% confidence level: $z^* = 2.05$.

$$4 \ge 2.05 * 18 / \sqrt{n} \rightarrow n \ge (2.05 * 18/4)^2 = 85.1$$

We know that the critical value associated with the 96% confidence level: $z^* = 2.05$.

$$4 \ge 2.05 * 18 / \sqrt{n} \rightarrow n \ge (2.05 * 18/4)^2 = 85.1$$

The minimum number of children required to attain the desired margin of error is 85.1. Since we can't sample 0.1 of a child, we must sample at least 86 children (round up, since rounding down to 85 would yield a slightly larger margin of error than desired).

OpenIntro Statistics, 2nd Edition

		Decision	
		fail to reject H_0	reject H_0
Truth	H ₀ true		
	H_A true		

		Decision	
		fail to reject H_0	reject H_0
Truth	H ₀ true		Type 1 Error, α
	H_A true		

 Type 1 error is rejecting H₀ when you shouldn't have, and the probability of doing so is α (significance level)

		Decision	
		fail to reject H_0	reject H_0
Truth	H ₀ true		Type 1 Error, α
	H_A true	Type 2 Error, β	

- Type 1 error is rejecting H₀ when you shouldn't have, and the probability of doing so is α (significance level)
- Type 2 error is failing to reject H₀ when you should have, and the probability of doing so is β (a little more complicated to calculate)

		Decision	
		fail to reject H_0	reject H_0
T	H ₀ true	$1 - \alpha$	Type 1 Error, α
Iruth	H_A true	Type 2 Error, β	

- Type 1 error is rejecting H₀ when you shouldn't have, and the probability of doing so is α (significance level)
- Type 2 error is failing to reject H₀ when you should have, and the probability of doing so is β (a little more complicated to calculate)
- *Power* of a test is the probability of correctly rejecting H_0 , and the probability of doing so is 1β

		Decision	
		fail to reject H_0	reject H_0
T	H ₀ true	$1 - \alpha$	Type 1 Error, α
Iruth	H_A true	Type 2 Error, β	Power, $1 - \beta$

- Type 1 error is rejecting H₀ when you shouldn't have, and the probability of doing so is α (significance level)
- Type 2 error is failing to reject H₀ when you should have, and the probability of doing so is β (a little more complicated to calculate)
- *Power* of a test is the probability of correctly rejecting H_0 , and the probability of doing so is 1β
- In hypothesis testing, we want to keep α and β low, but there are inherent trade-offs.

Type 2 error rate

If the alternative hypothesis is actually true, what is the chance that we make a Type 2 Error, i.e. we fail to reject the null hypothesis even when we should reject it?

- The answer is not obvious.
- If the true population average is very close to the null hypothesis value, it will be difficult to detect a difference (and reject H₀).
- If the true population average is very different from the null hypothesis value, it will be easier to detect a difference.
- Clearly, β depends on the *effect size* (δ)

Example - Blood Pressure

Blood pressure oscillates with the beating of the heart, and the systolic pressure is defined as the peak pressure when a person is at rest. The average systolic blood pressure for people in the U.S. is about 130 mmHg with a standard deviation of about 25 mmHg.

We are interested in finding out if the average blood pressure of employees at a certain company is <u>greater</u> than the national average, so we collect a random sample of 100 employees and measure their systolic blood pressure. What are the hypotheses?

Example - Blood Pressure

Blood pressure oscillates with the beating of the heart, and the systolic pressure is defined as the peak pressure when a person is at rest. The average systolic blood pressure for people in the U.S. is about 130 mmHg with a standard deviation of about 25 mmHg.

We are interested in finding out if the average blood pressure of employees at a certain company is <u>greater</u> than the national average, so we collect a random sample of 100 employees and measure their systolic blood pressure. What are the hypotheses?

 $H_0: \mu = 130$ $H_A: \mu > 130$

Example - Blood Pressure

Blood pressure oscillates with the beating of the heart, and the systolic pressure is defined as the peak pressure when a person is at rest. The average systolic blood pressure for people in the U.S. is about 130 mmHg with a standard deviation of about 25 mmHg.

We are interested in finding out if the average blood pressure of employees at a certain company is <u>greater</u> than the national average, so we collect a random sample of 100 employees and measure their systolic blood pressure. What are the hypotheses?

 $H_0: \mu = 130$ $H_A: \mu > 130$

We'll start with a very specific question – "What is the power of this hypothesis test to correctly detect an <u>increase</u> of 2 mmHg in average blood pressure?"
The preceding question can be rephrased as "How likely is it that this test will reject H_0 when the true average systolic blood pressure for employees at this company is 132 mmHg?"

The preceding question can be rephrased as "How likely is it that this test will reject H_0 when the true average systolic blood pressure for employees at this company is 132 mmHg?"

Hint: Break this down intro two simpler problems

The preceding question can be rephrased as "How likely is it that this test will reject H_0 when the true average systolic blood pressure for employees at this company is 132 mmHg?"

Hint: Break this down intro two simpler problems

1. Problem 1: Which values of \bar{x} represent sufficient evidence to reject H_0 ?

The preceding question can be rephrased as "How likely is it that this test will reject H_0 when the true average systolic blood pressure for employees at this company is 132 mmHg?"

Hint: Break this down intro two simpler problems

- 1. Problem 1: Which values of \bar{x} represent sufficient evidence to reject H_0 ?
- 2. Problem 2: What is the probability that we would reject H_0 if \bar{x} had come from $N\left(mean = 132, SE = \frac{25}{\sqrt{100}} = 2.5\right)$, i.e. what is the probability that we can obtain such an \bar{x} from this distribution?

The preceding question can be rephrased as "How likely is it that this test will reject H_0 when the true average systolic blood pressure for employees at this company is 132 mmHg?"

Hint: Break this down intro two simpler problems

- 1. Problem 1: Which values of \bar{x} represent sufficient evidence to reject H_0 ?
- 2. Problem 2: What is the probability that we would reject H_0 if \bar{x} had come from $N\left(mean = 132, SE = \frac{25}{\sqrt{100}} = 2.5\right)$, i.e. what is the probability that we can obtain such an \bar{x} from this distribution?

Determine how power changes as sample size, standard deviation of the sample, α , and effect size increases.

Which values of \bar{x} represent sufficient evidence to reject H_0 ? (Remember $H_0: \mu = 130, H_A: \mu > 130$)

Which values of \bar{x} represent sufficient evidence to reject H_0 ? (Remember $H_0: \mu = 130, H_A: \mu > 130$)



Which values of \bar{x} represent sufficient evidence to reject H_0 ? (Remember $H_0: \mu = 130, H_A: \mu > 130$)



Any $\bar{x} > 134.125$ would be sufficient to reject H_0 at the 5% significance level.

What is the probability that we would reject H_0 if \bar{x} did come from N(mean = 132, SE = 2.5).

What is the probability that we would reject H_0 if \bar{x} did come from N(mean = 132, SE = 2.5).

This is the same as finding the area above $\bar{x} = 134.125$ if \bar{x} came from N(132, 2.5).

What is the probability that we would reject H_0 if \bar{x} did come from N(mean = 132, SE = 2.5).

This is the same as finding the area above $\bar{x} = 134.125$ if \bar{x} came from N(132, 2.5).



What is the probability that we would reject H_0 if \bar{x} did come from N(mean = 132, SE = 2.5).

This is the same as finding the area above $\bar{x} = 134.125$ if \bar{x} came from N(132, 2.5).



The probability of rejecting H_0 : $\mu = 130$, if the true average systolic blood pressure of employees at this company is 132 mmHg, is 0.1977 which is the power of this test. Therefore, $\beta = 0.8023$ for this test.













There are several ways to increase power (and hence decrease type 2 error rate):

There are several ways to increase power (and hence decrease type 2 error rate):

1. Increase the sample size.

There are several ways to increase power (and hence decrease type 2 error rate):

- 1. Increase the sample size.
- 2. Decrease the standard deviation of the sample, which essentially has the same effect as increasing the sample size (it will decrease the standard error). With a smaller *s* we have a better chance of distinguishing the null value from the observed point estimate. This is difficult to ensure but cautious measurement process and limiting the population so that it is more homogenous may help.

There are several ways to increase power (and hence decrease type 2 error rate):

- 1. Increase the sample size.
- 2. Decrease the standard deviation of the sample, which essentially has the same effect as increasing the sample size (it will decrease the standard error). With a smaller *s* we have a better chance of distinguishing the null value from the observed point estimate. This is difficult to ensure but cautious measurement process and limiting the population so that it is more homogenous may help.
- 3. Increase α , which will make it more likely to reject H_0 (but note that this has the side effect of increasing the Type 1 error rate).

There are several ways to increase power (and hence decrease type 2 error rate):

- 1. Increase the sample size.
- 2. Decrease the standard deviation of the sample, which essentially has the same effect as increasing the sample size (it will decrease the standard error). With a smaller *s* we have a better chance of distinguishing the null value from the observed point estimate. This is difficult to ensure but cautious measurement process and limiting the population so that it is more homogenous may help.
- 3. Increase α , which will make it more likely to reject H_0 (but note that this has the side effect of increasing the Type 1 error rate).
- 4. Consider a larger effect size. If the true mean of the population is in the alternative hypothesis but close to the null value, it will be harder to detect a difference.

OpenIntro Statistics, 2nd Edition

Recap - Calculating Power

- Begin by picking a meaningful effect size δ and a significance level α
- Calculate the range of values for the point estimate beyond which you would reject H_0 at the chosen α level.
- Calculate the probability of observing a value from preceding step if the sample was derived from a population where $\bar{x} \sim N(\mu_{H_0} + \delta, SE)$

Going back to the blood pressure example, how large a sample would you need if you wanted 90% power to detect a 4 mmHg increase in average blood pressure for the hypothesis that the population average is greater than 130 mmHg at $\alpha = 0.05$?

Going back to the blood pressure example, how large a sample would you need if you wanted 90% power to detect a 4 mmHg increase in average blood pressure for the hypothesis that the population average is greater than 130 mmHg at $\alpha = 0.05$?

Given: $H_0: \mu = 130, H_A: \mu > 130, \alpha = 0.05, \beta = 0.10, \sigma = 25, \delta = 4$

Going back to the blood pressure example, how large a sample would you need if you wanted 90% power to detect a 4 mmHg increase in average blood pressure for the hypothesis that the population average is greater than 130 mmHg at $\alpha = 0.05$?

Given: $H_0: \mu = 130, H_A: \mu > 130, \alpha = 0.05, \beta = 0.10, \sigma = 25, \delta = 4$

Step 1: Determine the cutoff – in order to reject H_0 at $\alpha = 0.05$, we need a sample mean that will yield a Z score of at least 1.65.

$$\bar{x} > 130 + 1.65 \frac{25}{\sqrt{n}}$$

Going back to the blood pressure example, how large a sample would you need if you wanted 90% power to detect a 4 mmHg increase in average blood pressure for the hypothesis that the population average is greater than 130 mmHg at $\alpha = 0.05$?

Given:
$$H_0: \mu = 130$$
, $H_A: \mu > 130$, $\alpha = 0.05$, $\beta = 0.10$, $\sigma = 25$, $\delta = 4$

Step 1: Determine the cutoff – in order to reject H_0 at $\alpha = 0.05$, we need a sample mean that will yield a Z score of at least 1.65.

$$\bar{x} > 130 + 1.65 \frac{25}{\sqrt{n}}$$

Step 2: Set the probability of obtaining the above \bar{x} if the true population is centered at 130 + 4 = 134 to the desired power, and solve for *n*.

$$P\left(\bar{x} > 130 + 1.65\frac{25}{\sqrt{n}}\right) = 0.9$$
$$P\left(Z > \frac{\left(130 + 1.65\frac{25}{\sqrt{n}}\right) - 134}{\frac{25}{\sqrt{n}}}\right) = P\left(Z > 1.65 - 4\frac{\sqrt{n}}{25}\right) = 0.9$$







You can either directly solve for n, or use computation to calculate power for various n and determine the sample size that yields the desired power:



For n = 336, power = 0.9002, therefore we need 336 subjects in our sample to achieve the desired level of power for the given circumstance.

OpenIntro Statistics, 2nd Edition

- Variability in estimates
- 2 Confidence intervals
- 3 Hypothesis testing
- Examining the Central Limit Theorem
- Inference for other estimators
- Sample size and power



All else held equal, will the p-value be lower if n = 100 or n = 10,000?

- (a) n = 100
- (b) n = 10,000

All else held equal, will the p-value be lower if n = 100 or n = 10,000?

(a) n = 100

(b) n = 10,000
(a) n = 100

(b) n = 10,000

(a) n = 100

(b) n = 10,000

$$Z_{n=100} = \frac{50 - 49.5}{\frac{2}{\sqrt{100}}}$$

(a) n = 100

(b) n = 10,000

$$Z_{n=100} = \frac{50 - 49.5}{\frac{2}{\sqrt{100}}} = \frac{50 - 49.5}{\frac{2}{10}} = \frac{0.5}{0.2} = 2.5, \quad p\text{-value} = 0.0062$$

(a) n = 100

(b) n = 10,000

$$Z_{n=100} = \frac{50 - 49.5}{\frac{2}{\sqrt{100}}} = \frac{50 - 49.5}{\frac{2}{10}} = \frac{0.5}{0.2} = 2.5, \quad p\text{-value} = 0.0062$$
$$Z_{n=10000} = \frac{50 - 49.5}{\frac{2}{\sqrt{10000}}}$$

(a) n = 100

(b) n = 10,000

$$Z_{n=100} = \frac{50 - 49.5}{\frac{2}{\sqrt{100}}} = \frac{50 - 49.5}{\frac{2}{10}} = \frac{0.5}{0.2} = 2.5, \quad p\text{-value} = 0.0062$$
$$Z_{n=10000} = \frac{50 - 49.5}{\frac{2}{\sqrt{10000}}} = \frac{50 - 49.5}{\frac{2}{100}} = \frac{0.5}{0.02} = 25, \quad p\text{-value} \approx 0$$

(a) n = 100

(b) n = 10,000

Suppose $\bar{x} = 50$, s = 2, $H_0 : \mu = 49.5$, and $H_A : \mu \ge 49.5$.

$$Z_{n=100} = \frac{50 - 49.5}{\frac{2}{\sqrt{100}}} = \frac{50 - 49.5}{\frac{2}{10}} = \frac{0.5}{0.2} = 2.5, \quad p\text{-value} = 0.0062$$
$$Z_{n=10000} = \frac{50 - 49.5}{\frac{2}{\sqrt{10000}}} = \frac{50 - 49.5}{\frac{2}{100}} = \frac{0.5}{0.02} = 25, \quad p\text{-value} \approx 0$$

As *n* increases - $SE \downarrow$, $Z \uparrow$, *p*-value \downarrow

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30			
<i>n</i> = 5000			

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45		
<i>n</i> = 5000			

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45		
<i>n</i> = 5000	p - value = 0.04		

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45	p - value = 0.39	
<i>n</i> = 5000	p - value = 0.04		

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45	p - value = 0.39	
<i>n</i> = 5000	p - value = 0.04	p - value = 0.0002	

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45	p - value = 0.39	p - value = 0.29
<i>n</i> = 5000	p - value = 0.04	p - value = 0.0002	

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45	p - value = 0.39	p - value = 0.29
<i>n</i> = 5000	p - value = 0.04	p - value = 0.0002	$p - value \approx 0$

\bar{x}	10.05	10.1	10.2
<i>n</i> = 30	p - value = 0.45	p - value = 0.39	p - value = 0.29
<i>n</i> = 5000	p - value = 0.04	p - value = 0.0002	$p - value \approx 0$

When n is large, even small deviations from the null (small effect sizes), which may be considered practically insignificant, can yield statistically significant results.

Statistical vs. practical significance

- Real differences between the point estimate and null value are easier to detect with larger samples.
- However, very large samples will result in statistical significance even for tiny differences between the sample mean and the null value (*effect size*), even when the difference is not practically significant.
- This is especially important to research: if we conduct a study, we want to focus on finding meaningful results (we want observed differences to be real, but also large enough to matter).
- The role of a statistician is not just in the analysis of data, but also in planning and design of a study.

"To call in the statistician after the experiment is done may be no more than asking him to perform a post-mortem examination: he may be able to say what the experiment died of." – R.A. Fisher