

Chapter 5: Foundations for inference

OpenIntro Statistics, 4th Edition

Slides developed by Mine Çetinkaya-Rundel of OpenIntro.

The slides may be copied, edited, and/or shared via the CC BY-SA license.

Some images may be included under fair use guidelines (educational purposes).

Point estimates and sampling variability

Point estimates and error

- We are often interested in *population parameters*.
- Complete populations are difficult to collect data on, so we use *sample statistics* as *point estimates* for the unknown population parameters of interest.
- *Error* in the estimate = difference between population parameter and sample statistic
- *Bias* is systematic tendency to over- or under-estimate the true population parameter.
- *Sampling error* describes how much an estimate will tend to vary from one sample to the next.
- Much of statistics is focused on understanding and quantifying sampling error, and *sample size* is helpful for quantifying this error.

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.

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.

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% \pm 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% \pm 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.

Suppose the proportion of American adults who support the expansion of solar energy is $p = 0.88$, which is our parameter of interest. Is a randomly selected American adult more or less likely to support the expansion of solar energy?

More likely.

Suppose that you don't have access to the populaion of all American adults, which is a quite likely scenario. In order to estimate the proportion of American adults who support solar power expansion, you might sample from the population and use your sample proportion as the best guess for the unknown population proportion.

- Sample, with replacement, 1000 American adults from the population, and record whether they support solar power or not expansion.
- Find the sample proportion.
- Plot the distribution of the sample proportions obtained by members of the lass.

1. Create a set of 250 million entries, where 88\% of
them are "support" and 12\% are "not".

```
pop_size <- 2500000000  
possible_entries <- c(rep("support", 0.88 * pop_size),  
                      rep("not", 0.12 * pop_size))
```

2. Sample 1000 entries without replacement.

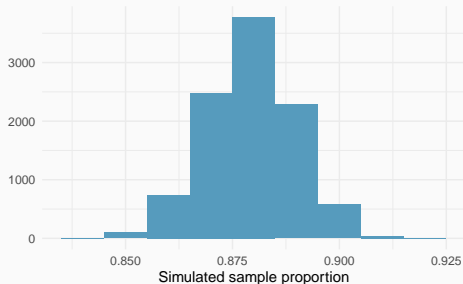
```
sampled_entries <- sample(possible_entries, size = 1000)
```

3. Compute \hat{p} : count the number that are "support",
then divide by # the sample size.

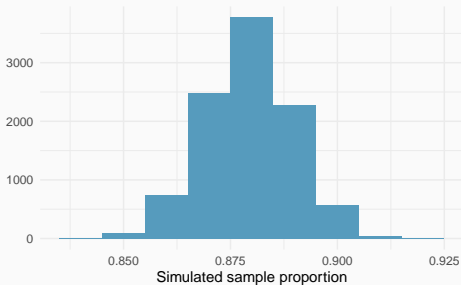
```
sum(sampled_entries == "support") / 1000
```

Sampling distribution

Suppose you were to repeat this process many times. The distribution of \hat{p} s you 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 proportion?



Approximately 0.88, the true population proportion.

Sampling distributions are never observed

- In real-world applications, we never actually observe the sampling distribution, yet it is useful to always think of a point estimate as coming from such a hypothetical distribution.
- Understanding the sampling distribution will help us characterize and make sense of the point estimates that we do observe.

Central Limit Theorem

Central limit theorem

Sample proportions will be nearly normally distributed with mean equal to the population mean, p , and standard error equal to

$$\sqrt{\frac{p(1-p)}{n}}.$$

$$\hat{p} \sim N\left(\text{mean} = p, SE = \sqrt{\frac{p(1-p)}{n}}\right)$$

- It wasn't a coincidence that the sampling distribution we saw earlier was symmetric, and centered at the true population population.
- We won't go through a detailed proof of why $SE = \sqrt{\frac{p(1-p)}{n}}$, but note that as n increases SE decreases.
 - As n increases samples will yield more consistent \hat{p} s, i.e. variability among \hat{p} s will 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.
2. *Sample size*: There should be at least 10 expected successes and 10 expected failures in the observed sample. This is difficult to verify if you don't know the population proportion (or can't assume a value for it). In those cases we look for the number of observed successes and failures to be at least 10.

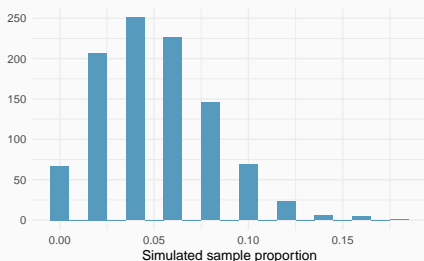
When p is unknown

- The CLT states $SE = \sqrt{\frac{p(1-p)}{n}}$, with the condition that np and $n(1-p)$ are at least 10, however we often don't know the value of p , the population proportion
- In these cases we substitute \hat{p} for p

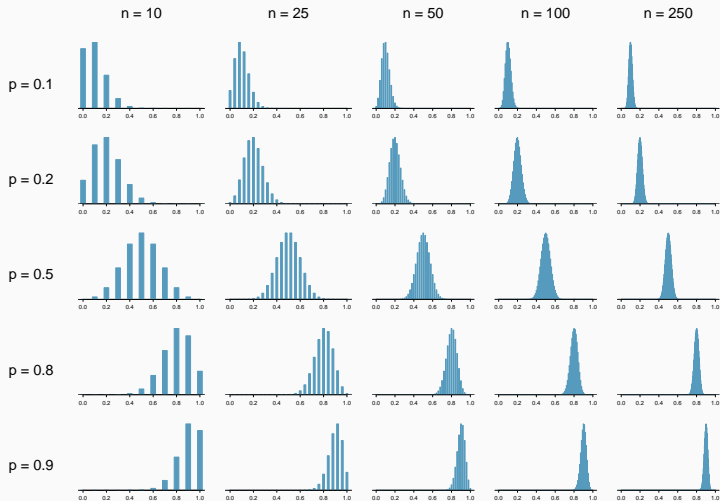
When p is low

Suppose we have a population where the true population proportion is $p = 0.05$, and we take random samples of size $n = 50$ from this population. We calculate the sample proportion in each sample and plot these proportions. Would you expect this distribution to be nearly normal? Why, or why not?

*No, the success-failure condition is not met ($50 * 0.05 = 2.5$), hence we would not expect the sampling distribution to be nearly normal.*



What happens when np and/or $n(1 - p) \downarrow 10$?



When the conditions are not met...

- When either np or $n(1 - p)$ is small, the distribution is more discrete.
- When np or $n(1 - p) < 10$, the distribution is more skewed.
- The larger both np and $n(1 - p)$, the more normal the distribution.
- When np and $n(1 - p)$ are both very large, the discreteness of the distribution is hardly evident, and the distribution looks much more like a normal distribution.

Extending the framework for other statistics

- The strategy of using a sample statistic to estimate a parameter is quite common, and it's a strategy that we can apply to other statistics besides a proportion.
 - Take a random sample of students at a college and ask them how many extracurricular activities they are involved in to estimate the average number of extra curricular activities all students in this college are interested in.
- The principles and general ideas are from this chapter apply to other parameters as well, even if the details change a little.

Confidence intervals for a proportion

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.

Facebook's categorization of user interests

Most commercial websites (e.g. social media platforms, news outlets, online retailers) collect a data about their users' behaviors and use these data to deliver targeted content, recommendations, and ads. To understand whether Americans think their lives line up with how the algorithm-driven classification systems categorizes them, Pew Research asked a representative sample of 850 American Facebook users how accurately they feel the list of categories Facebook has listed for them on the page of their supposed interests actually represents them and their interests. 67% of the respondents said that the listed categories were accurate. Estimate the true proportion of American Facebook users who think the Facebook categorizes their interests accurately.

<https://www.pewinternet.org/2019/01/16/facebook-algorithms-and-personal-data/>

Facebook's categorization of user interests

$$\hat{p} = 0.67 \quad n = 850$$

The approximate 95% confidence interval is defined as

$$\text{point estimate} \pm 1.96 \times SE$$

$$SE = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{0.67 \times 0.33}{850}} \approx 0.016$$

$$\begin{aligned}\hat{p} \pm 1.96 \times SE &= 0.67 \pm 1.96 \times 0.016 \\ &= (0.67 - 0.03, 0.67 + 0.03) \\ &= (0.64, 0.70)\end{aligned}$$

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

We are 95% confident that

- (a) 64% to 67% of American Facebook users in this sample think Facebook categorizes their interests accurately.
- (b) *64% to 67% of all American Facebook users think Facebook categorizes their interests accurately*
- (c) there is a 64% to 67% chance that a randomly chosen American Facebook user's interests are categorized accurately.
- (d) there is a 64% to 67% chance that 95% of American Facebook users' interests are categorized accurately.

What does 95% confident mean?

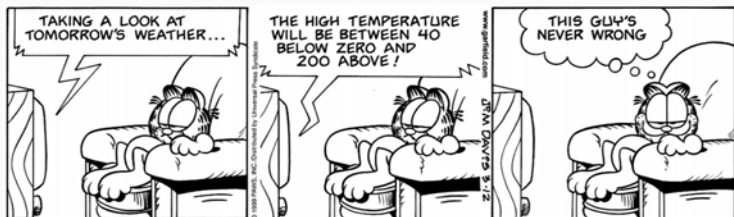
- Suppose we took many samples and built a confidence interval from each sample using the equation $\text{point estimate} \pm 1.96 \times SE$.
- Then about 95% of those intervals would contain the true population proportion (p).

Width of an 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.

Changing the confidence level

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

- In a confidence interval, $z^{\star} \times SE$ 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^{\star} in the above formula.
- Commonly used confidence levels in practice are 90%, 95%, 98%, and 99%.
- For a 95% confidence interval, $z^{\star} = 1.96$.
- However, using the standard normal (z) distribution, it is possible to find the appropriate z^{\star} for any confidence level.

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

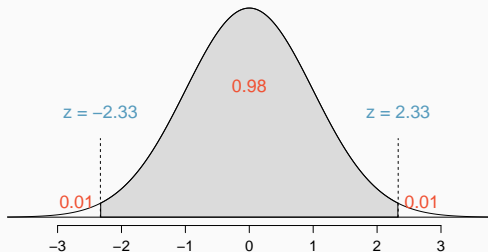
(a) $Z = 2.05$

(d) $Z = -2.33$

(b) $Z = 1.96$

(e) $Z = -1.65$

(c) $Z = 2.33$



Interpreting confidence intervals

Confidence intervals are ...

- always about the population
- are not probability statements
- only about population parameters, not individual observations
- only reliable if the sample statistic they're based on is an unbiased estimator of the population parameter

Hypothesis testing for a proportion

Remember when...

Gender discrimination experiment:

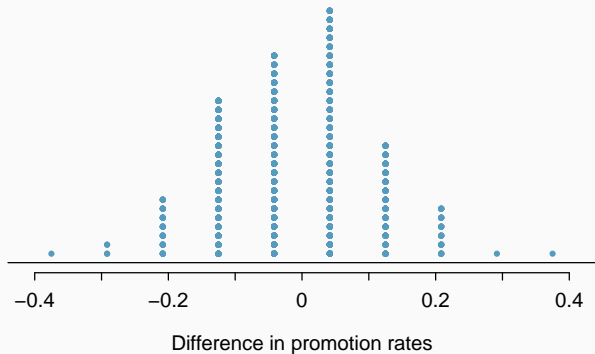
		<i>Promotion</i>		Total
		Promoted	Not Promoted	
<i>Gender</i>	Male	21	3	24
	Female	14	10	24
	Total	35	13	48

$$\hat{p}_{males} = 21/24 \approx 0.88 \text{ and } \hat{p}_{females} = 14/24 \approx 0.58$$

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



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.

Recap: hypothesis testing framework

- 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.

Testing hypotheses using confidence intervals

Earlier we calculated a 95% confidence interval for the proportion of American Facebook users who think Facebook categorizes their interests accurately as 64% to 67%. Based on this confidence interval, do the data support the hypothesis that majority of American Facebook users think Facebook categorizes their interests accurately.

- The associated hypotheses are:
 $H_0: p = 0.50$: 50% of American Facebook users think Facebook categorizes their interests accurately
 $H_A: p > 0.50$: More than 50% of American Facebook users think Facebook categorizes their interests accurately
- Null value is not included in the interval \rightarrow reject the null hypothesis.
- This is a quick-and-dirty approach for hypothesis testing, but it doesn't tell us the likelihood of certain outcomes under the null hypothesis (p-value).

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 errors (cont.)

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	✓	
	H_A true		✓

Decision errors (cont.)

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	✓	Type 1 Error
	H_A true		✓

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

Decision errors (cont.)

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	✓	Type 1 Error
	H_A true	Type 2 Error	✓

- 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 errors (cont.)

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	✓	Type 1 Error
	H_A true	Type 2 Error	✓

- 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.

Hypothesis Test as a trial

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 circumstances?

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

Hypothesis Test as a trial

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 circumstances?

- 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?

Hypothesis Test as a trial

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 circumstances?

- 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

Type 1 error rate

- 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.

$$P(\text{Type 1 error} \mid H_0 \text{ true}) = \alpha$$

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

Facebook interest categories

The same survey asked the 850 respondents how comfortable they are with Facebook creating a list of categories for them. 41% of the respondents said they are comfortable. Do these data provide convincing evidence that the proportion of American Facebook users are comfortable with Facebook creating a list of interest categories for them is different than 50%?

<https://www.pewinternet.org/2019/01/16/facebook-algorithms-and-personal-data/>

Setting the hypotheses

- The *parameter of interest* is the proportion of all American Facebook users who are comfortable with Facebook creating categories of interests for them.
- There may be two explanations why our sample proportion is lower than 0.50 (minority).
 - The true population proportion is different than 0.50.
 - The true population mean is 0.50, and the difference between the true population proportion and the sample proportion is simply due to natural sampling variability.

Setting the hypotheses

- We start with the assumption that 50% of American Facebook users are comfortable with Facebook creating categories of interests for them

$$H_0 : p = 0.50$$

- We test the claim that the proportion of American Facebook users who are comfortable with Facebook creating categories of interests for them is different than 50%

$$H_A : p \neq 0.50$$

Facebook interest categories - conditions

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

- (a) Respondents in the sample should be independent of each other with respect to whether or not they feel comfortable with their interests being categorized by Facebook.
- (b) Sampling should have been done randomly.
- (c) The sample size should be less than 10% of the population of all American Facebook users.
- (d) *There should be at least 30 respondents in the sample.*
- (e) There should be at least 10 expected successes and 10 expected failure.

Test statistic

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*.

$$\hat{p} \sim N\left(\mu = 0.50, SE = \sqrt{\frac{0.50 \times 0.50}{850}}\right)$$

$$Z = \frac{0.41 - 0.50}{0.0171} = -5.26$$

Test statistic

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*.

$$\hat{p} \sim N\left(\mu = 0.50, SE = \sqrt{\frac{0.50 \times 0.50}{850}}\right)$$

$$Z = \frac{0.41 - 0.50}{0.0171} = -5.26$$

The sample proportion is 5.26 standard errors away from the hypothesized value. Is this considered unusually low? 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.
- 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 H_0* .
- 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* .

Facebook interest categories - Making a decision

- $p\text{-value} < 0.0001$
 - If 50% of all American Facebook users are comfortable with Facebook creating these interest categories, there is less than a 0.01% chance of observing a random sample of 850 American Facebook users where 41% or fewer or 59% or higher feel comfortable with it.
 - This is a pretty low probability for us to think that the observed sample proportion, or something more extreme, is likely to happen simply by chance.
- Since $p\text{-value}$ is *low* (lower than 5%) we *reject H_0* .
- The data provide convincing evidence that the proportion of American Facebook users who are comfortable with Facebook creating a list of interest categories for them is different than 50%.

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

One vs. two sided hypothesis tests

- In two sided hypothesis tests we are interested in whether p is either above or below some null value p_0 : $H_A : p \neq p_0$.
- In one sided hypothesis tests we are interested in p differing from the null value p_0 in one direction (and not the other):
 - If there is only value in detecting if population parameter is less than p_0 , then $H_A : p < p_0$.
 - If there is only value in detecting if population parameter is greater than p_0 , then $H_A : p > p_0$.
- Two-sided tests are often more appropriate as we often want to detect if the data goes clearly in the opposite direction of our alternative hypothesis as well.