

## 8. Hypothesis testing: two groups

LPO.8800: Statistical Methods in Education Research

Sean P. Corcoran

### Hypothesis testing thus far

- Hypothesis tests for a population mean  $\mu$  ( $\sigma$  known or unknown)
- Hypothesis tests for a population proportion  $\pi$
- Null and alternative hypotheses ( $H_0$  and  $H_1$ )
- One vs. two-sided alternative hypotheses
- Significance levels ( $\alpha$ ) and  $p$ -values
- Type I and Type II errors, Power
- Practical significance and effect size
- $p$ -screening and  $p$ -hacking
- Multiple hypothesis testing

# Statistical tests for comparing two groups

Hypothesis tests are frequently used to make inferences about how two population parameters compare:

- Do female executives earn less on average than males?
- Do 4th graders in an experimental reading program perform differently on standardized reading tests than 4th graders not in the program?
- Are women more likely to vote for Democratic candidates than men?
- Do subjects participating in a 6-week weight loss program lose more weight over time than those who do not participate in the program?
- Has obesity among children aged 10-12 increased between 2010 and 2020?
- Were COVID infection rates higher in counties without a mask mandate than counties with them?

## Statistical tests for comparing two groups

We can think of these as **bivariate analyses** involving two variables:

- Group identifier: a binary *explanatory variable*
- Outcome: the *response*

## Statistical tests for comparing two groups

Each of the above examples is a comparison of two parameters. For example:

- A comparison of means across groups ( $\mu_1$  and  $\mu_2$ )
- A comparison of proportions across groups ( $\pi_1$  and  $\pi_2$ )
- A comparison of means or proportions *over time*
- A comparison of means of the *same group* pre- and post-treatment (“within subject”)

When making comparisons of two parameters, we actually construct a test for the *difference* in those parameters (e.g.,  $\mu_2 - \mu_1$ , or  $\pi_2 - \pi_1$ ).

## Statistical tests for comparing two groups

The steps for conducting a test comparing two groups are the same as those for the test of a single parameter. The most common null hypothesis is that there is *no difference* between the two population means:

$$H_0 : \mu_1 = \mu_2$$

Equivalently,

$$H_0 : \mu_2 - \mu_1 = 0$$

Note: it doesn't matter which mean you subtract from the other, as long as you keep track and are consistent in your ordering throughout the test.

## Statistical tests for comparing two groups

The alternative hypothesis  $H_1$  is that there *is* a difference (a two-sided alternative) or that one mean is greater than the other (a one-sided alternative):

Two-sided alternative:

$$H_1 : \mu_2 - \mu_1 \neq 0$$

One-sided alternatives:

$$H_1 : \mu_2 - \mu_1 > 0$$

$$H_1 : \mu_2 - \mu_1 < 0$$

### Hypothesis test steps

- 1 Determine  $H_0$  and  $H_1$ .
- 2 The estimator for a difference in population means is the difference in sample means: e.g.,  $\bar{x}_2 - \bar{x}_1$ . Determine the sampling distribution of your test statistic under  $H_0$ . (This requires knowing the standard error of  $\bar{x}_2 - \bar{x}_1$ ).
- 3 Determine the probability of obtaining your observed test statistic if  $H_0$  is true (the  $p$ -value), and draw a conclusion.

The standard error of  $\bar{x}_2 - \bar{x}_1$  depends on how the two samples were drawn (see next slide).

## Independent vs. dependent samples

The design of a study—and in particular, whether the two samples being compared are independent or dependent—is important to statistical testing, because standard errors depend on how the samples were drawn.

- Samples are **dependent** when there is a natural matching between subjects in each sample. Examples include repeated measures on the *same* subjects (a **longitudinal study**), siblings, marriage partners, etc.
- Samples are **independent** when there is no such matching (e.g., random draws from two populations). Selection of subjects into one sample has no effect on the selection of subjects into the second.

With dependent samples, pairs of outcomes do not represent independent draws from a population—they are likely to be correlated.

## Experimental vs. observational group assignment

- In **experimental** designs, subjects are *randomly* assigned to groups (e.g., treatment and control).
- In **observational** designs, subjects are in naturally-occurring groups that they may or may not have control over (e.g., gender or political affiliation).

The distinction is important for causal inference and interpretation, though not necessarily for statistical inference (e.g., comparing population means).

## Sampling distribution of $\bar{x}_2 - \bar{x}_1$

Over repeated samples, the difference in two means drawn from **independent samples** will have a mean of  $\mu_2 - \mu_1$  (the difference in the *true* means—it is an *unbiased* estimator) and a standard error of:

$$se_{\bar{x}_2 - \bar{x}_1} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Notice this is a larger number than either of the individual standard errors alone (for  $\bar{x}_1$  or  $\bar{x}_2$ )! Also note this expression is:  $\sqrt{(se_1)^2 + (se_2)^2}$

Intuitively,  $(\bar{x}_2 - \bar{x}_1)$  tends to be further from  $(\mu_2 - \mu_1)$  than  $\bar{x}_1$  is from  $\mu_1$  or  $\bar{x}_2$  is from  $\mu_2$ .

## Sampling distribution of $\bar{x}_2 - \bar{x}_1$

If the samples are independent *and* the sample sizes  $n_1$  and  $n_2$  are sufficiently large, then we can also say the sampling distribution of  $\bar{x}_2 - \bar{x}_1$  (divided by its standard error) has an approximately normal distribution.

In practice,  $\sigma_1$  and  $\sigma_2$  are unknown, so the sample standard deviations are used in their place ( $s_1$  and  $s_2$ ):

$$se_{\bar{x}_2 - \bar{x}_1} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

We then refer to the *t*-distribution instead of the standard normal. (Which *df* to use is discussed momentarily).

## Confidence interval for $\mu_2 - \mu_1$

With this information, a  $(1 - \alpha)\%$  confidence interval for the difference in population means  $\mu_2 - \mu_1$  is:

$$(\bar{x}_2 - \bar{x}_1) \pm t_{\alpha/2}(se_{\bar{x}_2 - \bar{x}_1})$$

As before, a  $(1 - \alpha)\%$  confidence interval can be used to test two-sided alternative hypotheses with significance level  $\alpha$ .

## Hypothesis test for $\mu_2 - \mu_1$

Alternatively, a test statistic can be calculated for a specific hypothesis. For example:

$$H_0 : \mu_2 - \mu_1 = 0$$

$$H_1 : \mu_2 - \mu_1 \neq 0$$

Under the assumption that  $H_0$  is true, the test statistic is as follows, using the standard error formula given above:

$$t = \frac{(\bar{x}_2 - \bar{x}_1) - 0}{se_{\bar{x}_2 - \bar{x}_1}}$$

Use the  $t$  distribution to find the  $p$ -value associated with this  $t$ .

## Hypothesis test for $\mu_2 - \mu_1$

The degrees of freedom for the  $t$ -statistic in this case is complex (use the Satterthwaite approximation). However, if we can safely assume that the two distributions from which the samples are drawn have *equal variances* (**homoskedasticity**), the degrees of freedom will be  $df = n_1 + n_2 - 2$ .

Stata's default is to assume equal variances and use  $df = n_1 + n_2 - 2$ . Note this might not be the right assumption your case! Can use the unequal variance assumption.

## Satterthwaite approximation

The Satterthwaite approximation for the degrees of freedom in a two-sample independent  $t$ -test is

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{1}{n_1-1} \left( \frac{s_1^2}{n_1} \right)^2 + \frac{1}{n_2-1} \left( \frac{s_2^2}{n_2} \right)^2}$$



# Satterthwaite approximation

A few observations about the Satterthwaite approximation:

- When the variances are equal  $s_1^2 = s_2^2$  and the sample sizes are equal  $n_1 = n_2$  the Satterthwaite  $df = n_1 + n_2 - 2$
- Holding the variances (and total  $N$ ) constant, the more imbalanced are the two samples  $n_1$  and  $n_2$ , the smaller is the Satterthwaite  $df$
- Holding the sample sizes constant, the more different are the variances  $s_1^2$  and  $s_2^2$ , the smaller is the Satterthwaite  $df$

With large samples, statistical significance rarely hinges on this decision.

## Hypothesis test for $\mu_2 - \mu_1$

Note: if we can assume equal variances, then we can use a pooled estimate of the variance ( $s_p^2$ ), and the standard error for  $\bar{x}_2 - \bar{x}_1$  simplifies to:

$$se_{\bar{x}_2 - \bar{x}_1} = \sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

The pooled variance estimate is a weighted average of  $s_1^2$  and  $s_2^2$ :

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 - 1) + (n_2 - 1)}$$

Note:  $s_p^2$  is close—but not the same as—the combined variance (i.e., if you were to combine the data and treat it like one sample).

## Example 1

What is the impact of alcoholism during pregnancy on the IQ of infants? IQs of infants of 6 women with alcoholism (group 1) were compared with those of infants of 46 women without alcoholism (group 2).

$$n_1 = 6, \bar{x}_1 = 78, s_1^2 = 361$$

$$n_2 = 46, \bar{x}_2 = 99, s_2^2 = 256$$

$$H_0 : \mu_2 - \mu_1 = 0$$

$$H_1 : \mu_2 - \mu_1 \neq 0$$

## Example 1

The 95% confidence interval is

$$(\bar{x}_2 - \bar{x}_1) \pm t_{\alpha/2}(se_{\bar{x}_2 - \bar{x}_1})$$

Assuming equal variances, the  $t$ -statistic used here is  $t(50, 0.025) = 2.01$ , where  $df = n_1 + n_2 - 2$ . The standard error is:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 - 1) + (n_2 - 1)} = \frac{(5)361 + (45)256}{5 + 45} = 266.5$$

$$se_{\bar{x}_2 - \bar{x}_1} = \sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} = \sqrt{266.5 \left( \frac{1}{6} + \frac{1}{46} \right)} = 7.086$$

## Example 1

So:

$$(\bar{x}_2 - \bar{x}_1) \pm t_{\alpha/2}(se_{\bar{x}_2 - \bar{x}_1})$$

$$(99 - 78) \pm 2.01(7.086) = (6.757, 35.242)$$

The confidence interval does not contain zero so we can reject  $H_0$  at the  $\alpha = 0.05$  significance level.

## Example 1

Alternatively, can calculate the test statistic:

$$t = \frac{\bar{x}_2 - \bar{x}_1 - 0}{se_{\bar{x}_2 - \bar{x}_1}}$$

$$t = \frac{99 - 78 - 0}{7.086} = 2.964$$

The reference distribution for looking up the  $p$ -value is  $t(50)$ . In Stata: `display 2*ttail(50,2.964) = 0.0046`. Since  $p < \alpha$ , we reject  $H_0$ .

## Example 1: Using Stata `ttesti`

Can use the *t*-test calculator in Stata to obtain these results. Syntax below or use the drop-down menus. Default assumes equal variances, and *n*, *m*, *s* below refer to the two group sample sizes, means, and standard deviations.

`ttesti n1 m1 s1 n2 m2 s2`

```
. ttesti 6 78 19 46 99 16
```

Two-sample t test with equal variances

	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
x	6	78	7.756718	19	58.06072	97.93928
y	46	99	2.359071	16	94.24859	103.7514
combined	52	96.57692	2.430458	17.52628	91.69758	101.4563
diff		-21	7.085912		-35.23247	-6.767527

```
diff = mean(x) - mean(y)                                t = -2.9636
Ho: diff = 0                                             degrees of freedom = 50
Ha: diff < 0                                           Ha: diff != 0
Pr(T < t) = 0.0023                                     Pr(|T| > |t|) = 0.0046
Pr(T > t) = 0.9977
```

## Example 1: Using Stata `ttesti`

This is a case where the equal variance assumption makes a big difference:

`ttesti n1 m1 s1 n2 m2 s2, unequal`

```
. ttesti 6 78 19 46 99 16, unequal
```

Two-sample t test with unequal variances

	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
x	6	78	7.756718	19	58.06072	97.93928
y	46	99	2.359071	16	94.24859	103.7514
combined	52	96.57692	2.430458	17.52628	91.69758	101.4563
diff		-21	8.10752		-40.86902	-1.130985

```
diff = mean(x) - mean(y)                                t = -2.5902
Ho: diff = 0                                             Satterthwaite's degrees of freedom = 5.96208
Ha: diff < 0                                           Ha: diff != 0
Pr(T < t) = 0.0207                                     Pr(|T| > |t|) = 0.0414
Pr(T > t) = 0.9793
```

Bigger SE, smaller *t*-stat, fewer *df*, larger *p*, wider CI for diff.

## Example 1: Using Stata `ttesti`

Satterthwaite approximation for the degrees of freedom in this case:

$$df = \frac{\left(\frac{361}{6} + \frac{256}{46}\right)^2}{\frac{1}{6-1} \left(\frac{361}{6}\right)^2 + \frac{1}{46-1} \left(\frac{256}{46}\right)^2}$$
$$df = 5.96$$

This is substantially less than  $n_1 + n_2 - 2 = 50$ . Why? The sample sizes are very imbalanced, with one sample quite small ( $n_1 = 6$ ). Intuitively, with unequal variances, we have to use our data to estimate both—and one sample has only six observations.

## Example 1: Effect size

Is the difference in IQ *practically significant*? The Cohen's  $d$  measure is:

$$d = \frac{78 - 99}{17.5} = -1.2$$

The difference in IQ is 1.2 standard deviations in the (pooled) distribution of IQ: quite large! If you don't want to include the "treated" cases in the standard deviation calculation, could divide by 16 instead (the standard deviation for infants with mothers without alcoholism).

## Example 2: Using Stata

Use NELS to test the hypothesis that male and female college-bound high school graduates in the South have equal years of preparation in high school math (*unitmath*). Use  $\alpha=0.05$ . Let group 2 be males and group 1 be females.

$$H_0 : \mu_2 - \mu_1 = 0$$

$$H_1 : \mu_2 - \mu_1 \neq 0$$

Note: in the NELS, college-bound is *edexpect*  $\geq 2$  and the South is *region*  $= 3$ .

## Example 2: Using Stata

In Stata use `ttest varname, by(groupvar)`. The default assumes equal variances. Use the `unequal` option otherwise (this will affect the standard error calculation).

```
. ttest unitmath if edexpect>=2 & region==3, by(gender)
```

Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
male	57	3.938596	.0928778	.7012118	3.75254	4.124653
female	76	3.747105	.0783084	.6826767	3.591107	3.903104
combined	133	3.829173	.0602281	.6945841	3.710036	3.94831
diff		.1914912	.121017		-.0479093	.4308918

diff = mean(male) - mean(female) t = 1.5823  
Ho: diff = 0 degrees of freedom = 131

Ha: diff < 0  
Pr(T < t) = 0.9420

Ha: diff != 0  
Pr(|T| > |t|) = 0.1160

Ha: diff > 0  
Pr(T > t) = 0.0580

## Example 2: Using Stata

From the above Stata output:

$$\bar{x}_2 - \bar{x}_1 = 0.1915$$

$$se_{\bar{x}_2 - \bar{x}_1} = 0.1210$$

$$t = 1.582$$

$$p = 0.116$$

Conclusion: do not reject  $H_0$ , since  $p > \alpha$ .

## Example 2: Using Stata

Note: the above Stata output reports the combined standard deviation (0.6945841). The square of this (0.482447) is close, but is not the same as the pooled variance used to calculate the standard error of the difference:

$$s_p^2 = \frac{(57 - 1) * 0.7012^2 + (76 - 1) * 0.6827^2}{(57 - 1) + (76 - 1)} = 0.4770$$

Then:

$$se_{\bar{x}_2 - \bar{x}_1} = \sqrt{0.4770 \left( \frac{1}{57} + \frac{1}{76} \right)} = 0.121017$$

## Example 2: Using Stata

In practice, standard errors, test statistics and  $p$ -values will be similar whether assuming equal variances or not, if:

$n_1$  and  $n_2$  are similar or

$s_1^2$  and  $s_2^2$  are similar

## Example 2: Using Stata

Same example, using unequal option

```
. ttest unitmath if edexpect>=2 & region==3, by(gender) unequal
```

Two-sample t test with unequal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
male	57	3.938596	.0928778	.7012118	3.75254	4.124653
female	76	3.747105	.0783084	.6826767	3.591107	3.903104
combined	133	3.829173	.0602281	.6945841	3.710036	3.94831
diff		.1914912	.1214845		-.04906	.4320424

diff = mean(male) - mean(female)

t = 1.5763

Ho: diff = 0

Satterthwaite's degrees of freedom = 119.011

Ha: diff < 0

Ha: diff != 0

Ha: diff > 0

Pr(T < t) = 0.9412

Pr(|T| > |t|) = 0.1176

Pr(T > t) = 0.0588

Standard error in this case is  $\sqrt{(se_1)^2 + (se_2)^2}$ , not the calculation with  $s_p^2$



## Example 2: Effect size

Is the difference in math units *practically significant*? The Cohen's  $d$  measure is

$$d = \frac{3.94 - 3.75}{0.695} = 0.273$$

The male - female difference in math units is 0.19. This is about 0.273 standard deviations in the (pooled) distribution. Arguably a large effect, but would be worth benchmarking against other gaps.

## Tests for difference in proportions $\pi_2 - \pi_1$

What if the two parameters of interest are proportions? The procedure is the same:

- 1 Determine  $H_0$  and  $H_1$ .
- 2 The estimator for a difference in population means is the difference in sample proportions: e.g.,  $\hat{\pi}_2 - \hat{\pi}_1$ . Determine the sampling distribution of your test statistic under  $H_0$ .
- 3 Determine the probability of obtaining your observed test statistic if  $H_0$  is true (the  $p$ -value), and draw a conclusion.

## Sampling distribution of $\hat{\pi}_2 - \hat{\pi}_1$

Over repeated samples, the difference in two proportions drawn from independent samples will have a mean of  $\pi_2 - \pi_1$  (the *true* difference in means) and a standard error of:

$$se_{\hat{\pi}_2 - \hat{\pi}_1} = \sqrt{\frac{\hat{\pi}_1(1 - \hat{\pi}_1)}{n_1} + \frac{\hat{\pi}_2(1 - \hat{\pi}_2)}{n_2}}$$

If the samples are independent *and* the sample sizes  $n_1$  and  $n_2$  are sufficiently large, then we can also say the sampling distribution of  $\hat{\pi}_2 - \hat{\pi}_1$  has an approximately normal distribution. (I assume this below, in using  $z$  in the confidence interval). A typical rule of thumb is at least 10 successes and failures in each group:  $np > 10$  and  $n(1 - p) > 10$

## Confidence interval for $\pi_2 - \pi_1$

With this information, a  $(1 - \alpha)\%$  confidence interval for the difference in population proportions  $\pi_2 - \pi_1$  is:

$$(\hat{\pi}_2 - \hat{\pi}_1) \pm z_{\alpha/2}(se_{\hat{\pi}_2 - \hat{\pi}_1})$$

A  $(1 - \alpha)\%$  confidence interval can be used to test two-sided alternative hypotheses with significance level  $\alpha$ .

## Hypothesis test for $\pi_2 - \pi_1$

Alternatively, a test statistic can be calculated for a specific hypothesis.  
For example:

$$H_0 : \pi_2 - \pi_1 = 0$$

$$H_1 : \pi_2 - \pi_1 \neq 0$$

Under the assumption that  $H_0$  is true, the test statistic is as follows, using a modified standard error  $se_0$  (see next slide).

$$z = \frac{\hat{\pi}_2 - \hat{\pi}_1 - 0}{se_0}$$

## Hypothesis test for $\pi_2 - \pi_1$

Under  $H_0$  the two population proportions are the same. This directly implies equal variances, since the variance for proportions is  $\pi(1 - \pi)$ . If  $H_0$  is true, we can use a *pooled estimate* of  $\pi$  in the standard error formula that uses the combined samples rather than separate estimates of  $\pi_2$  and  $\pi_1$ . The standard error using the pooled estimate is:

$$se_0 = \sqrt{\hat{\pi}(1 - \hat{\pi}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

where  $\hat{\pi}$  is the overall proportion equal to one in the combined sample. This is used in the test statistic ( $z$ ) above.

## Example 3

An October 25, 2020 YouGov poll found 44% of 674 men and 36% of 826 women stated that they intend to vote for Donald Trump. Construct a 95% confidence interval for the difference in these two population proportions (men - women).

$$(\hat{\pi}_2 - \hat{\pi}_1) \pm z_{\alpha/2}(se_{\hat{\pi}_2 - \hat{\pi}_1})$$

$$(0.44 - 0.36) \pm 1.96(se_{\hat{\pi}_2 - \hat{\pi}_1})$$

$$se_{\hat{\pi}_2 - \hat{\pi}_1} = \sqrt{\frac{\hat{\pi}_1(1 - \hat{\pi}_1)}{n_1} + \frac{\hat{\pi}_2(1 - \hat{\pi}_2)}{n_2}} = \sqrt{\frac{0.36(0.64)}{826} + \frac{0.44(0.56)}{674}} = 0.0254$$

## Example 3

So:

$$(\hat{\pi}_2 - \hat{\pi}_1) \pm z_{\alpha/2}(se_{\hat{\pi}_2 - \hat{\pi}_1})$$

$$(0.44 - 0.36) \pm 1.96(0.0254)$$

$$0.08 \pm 1.96(0.0254) = (0.030, 0.130)$$

The confidence interval does not contain zero so we can conclude there is a statistically significant difference in support for Donald Trump between men and women.

## Example 3: using Stata `prtesti`

Can use the *t*-test calculator in Stata (for proportion) to obtain these results. Syntax below or use the drop-down menus. *n* and *p* refer to the two group sample sizes and proportions.

```
prtesti n1 p1 n2 p2
```

```
. prtesti 674 0.44 826 0.36
```

Two-sample test of proportions						x: Number of obs =	674
						y: Number of obs =	826
	Mean	Std. Err.	z	P> z	[95% Conf. Interval]		
x	.44	.0191201			.4025253	.4774747	
y	.36	.0167013			.327266	.392734	
diff	.08	.0253873			.0302419	.1297581	
	under Ho:	.0253853	3.15	0.002			
diff = prop(x) - prop(y)						z =	3.1514
Ho: diff = 0							
Ha: diff < 0			Ha: diff != 0		Ha: diff > 0		
Pr(Z < z) = 0.9992			Pr( Z  >  z ) = 0.0016		Pr(Z > z) = 0.0008		

## Example 4: Using Stata

Use NELS to test the hypothesis that male and female high school students differ in their propensity to binge drink. The variable *alcbinge* is equal to 1 if the student has ever binged on alcohol, and equal to 0 otherwise. Use  $\alpha=0.05$ . Let group 2 be males and group 1 be females.

$$H_0 : \pi_2 - \pi_1 = 0$$

$$H_1 : \pi_2 - \pi_1 \neq 0$$

## Example 4: Using Stata

In Stata use `prtest varname, by(groupvar)`. Note `prtest` for proportions.

```
. prtest alcbinge, by(gender)
```

Two-sample test of proportions

**male:** Number of obs = **227**

**female:** Number of obs = **273**

Group	Mean	Std. Err.	z	P> z	[95% Conf. Interval]
male	.2907489	.0301402			.2316752 .3498226
female	.1538462	.0218367			.111047 .1966453
diff	.1369027	.0372192			.0639544 .2098511
	under Ho:	.0369637	3.70	0.000	

diff = prop(male) - prop(female) z = 3.7037

Ho: diff = 0

Ha: diff < 0

Pr(Z < z) = 0.9999

Ha: diff != 0

Pr(|Z| > |z|) = 0.0002

Ha: diff > 0

Pr(Z > z) = 0.0001

## Example 4: Using Stata

From the above Stata output:

$$\hat{\pi}_2 - \hat{\pi}_1 = 0.1369$$

$$se_{\hat{\pi}_2 - \hat{\pi}_1} = 0.0369 \text{ under } H_0$$

$$t = 3.704$$

$$p = 0.0001$$

Conclusion: Reject  $H_0$ , since  $p < \alpha$ . There is a statistically significant difference in the propensity to binge drink between males and females. Note different SE used for the CI.

## Example 4: Using Stata

Note `ttest` results will differ from `prtest` for proportions—compare standard errors, for example. Better to use `prtest`.

```
. ttest alcbinge, by(gender)
```

Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
male	227	.2907489	.0302068	.4551114	.2312259	.3502719
female	273	.1538462	.0218768	.3614638	.1107768	.1969155
combined	500	.216	.0184219	.4119264	.1798059	.2521941
diff		.1369027	.0365263		.0651382	.2086673

diff = mean(male) - mean(female) t = 3.7481  
Ho: diff = 0 degrees of freedom = 498

Ha: diff < 0 Ha: diff != 0 Ha: diff > 0  
Pr(T < t) = 0.9999 Pr(|T| > |t|) = 0.0002 Pr(T > t) = 0.0001

## Paired sample *t*-test

- Suppose instead our samples are *dependent*, or paired.
- One advantage of a paired design is the ability to control for some external differences between the two samples.
- Independent samples will differ for a lot of idiosyncratic reasons (“noise”)
  - ▶ Paired samples allow you to “difference out” some of the noise that produces differences in independent sample means
  - ▶ One example: the *pre-post* design. Comparing the same individuals before and after some intervention eliminates many of the reasons samples differ

## Paired sample $t$ -test

With paired samples, we can conduct a test for the the average *within pair* difference, rather than the difference in two sample means. Here a null hypothesis of a zero difference in means is stated as:

$$H_0 : \mu_d = 0$$

$$H_1 : \mu_d \neq 0$$

where  $\mu_d$  is the mean within-pair difference, or mean “difference score.” A test of this hypothesis is simply the one-sample  $t$ -test for  $\bar{x}_d$ , the sample mean within-pair difference (calculated using each matched pair).

## Paired sample $t$ -test

The  $t$ -statistic for the paired sample  $t$ -test is:

$$t = \frac{\bar{x}_d - 0}{se_{\bar{x}_d}}$$

with  $n - 1$  degrees of freedom, where  $n$  is the number of matched pairs. The standard error of the mean difference is calculated as:

$$se_{\bar{x}_d} = s_d / \sqrt{n}$$

$s_d$  is the standard deviation of the paired differences.



## Example 5

In Everitt (1994), 17 girls treated for anorexia were weighed before and after treatment. Difference scores were calculated for each participant, with the following results:  $\bar{x}_d = 7.26$ ,  $s_d = 7.16$ . Test the null hypothesis that there was no change in weight.

$$H_0 : \mu_d = 0$$

$$H_1 : \mu_d \neq 0$$

## Example 5

The test statistic is:

$$t = \frac{7.26 - 0}{7.16/\sqrt{17}} = 4.17$$

With  $df=16$  the probability of obtaining by chance a  $t$ -statistic of 4.17 or larger is  $p < 0.01$ . With a significance level of  $\alpha = 0.05$  we reject  $H_0$ . The change in weight was statistically significant.

## Example 5 - confidence interval approach

Assuming a significance level of  $\alpha = 0.05$  we can alternatively construct a 95% confidence interval as:

$$\begin{aligned}\bar{x}_d \pm t(df)_{\alpha/2}(se_{\bar{x}_d}) \\ \bar{x}_d \pm t(16)_{0.025}(s_d/\sqrt{n}) \\ 7.26 \pm 2.12 \left( \frac{7.16}{\sqrt{17}} \right)\end{aligned}$$

(3.57, 10.95) - this interval does not contain zero, so we reject  $H_0$ . (The change in weight was statistically significant, or significantly different from zero).

## Example 6: Using Stata

Using the Agresti & Finlay dataset `anorexia.dta`, test for an effect of cognitive behavioral therapy on the weight of anorexia patients.  $n = 29$  subjects had `therapy==b`; *before* is the subject's weight before therapy and *after* is the subject's weight after therapy.

```
. ttest before=after if therapy=="b"
```

Paired t test						
Variable		Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]
before		29	82.68966	.8997857	4.845494	80.84653 84.53278
after		29	85.69655	1.550913	8.351924	82.51965 88.87345
diff		29	-3.006896	1.357155	7.308504	-5.786902 -.2268896

  

mean(diff) = mean(before - after)		t =	-2.2156
Ho: mean(diff) = 0		degrees of freedom =	28
Ha: mean(diff) < 0	Ha: mean(diff) != 0	Ha: mean(diff) > 0	
Pr(T < t) = 0.0175	Pr( T  >  t ) = 0.0350	Pr(T > t) = 0.9825	

## Example 6: Using Stata

Notice the change in syntax in the `ttest` command. From the above Stata output:

$$\bar{x}_d = -3.007$$

$$se_{\bar{x}_d} = 1.357$$

$$t = -2.216$$

$$p = 0.035$$

Conclusion: Reject  $H_0$ , since  $p < \alpha$ . The CBT therapy had a statistically significant effect on the subjects' weight.

## Stata `ttest` syntax recap

Stata `ttest` syntax:

- Test about a single population mean: `ttest varname==[#]`, where `[#]` is the population mean under  $H_0$ .
- Independent samples test comparing two population means: `ttest varname, by(group)`, where `group` is the group variable.
- Paired sample test: `ttest varname1==varname2`.
- Using Stata as a  $t$ -test calculator for a single population mean:  
`ttesti #obs #mean #sd #val [, level(#)]`
- Using Stata as a  $t$ -test calculator, independent samples test: `ttesti #obs1 #mean1 #sd1 #obs2 #mean2 #sd2 [, level(#)]`

# Stata prtest syntax recap

Stata prtest syntax:

- Test about a single population proportion: `prtest varname ==[#]`, where  $[\#]$  is the population proportion under  $H_0$ .
- Independent samples test comparing two population proportions: `prtest varname, by(group)`, where *group* is the group variable.
- Paired sample test (proportions): `prtest varname1==varname2`.
- Using Stata as a *t*-test calculator for a single population proportion: `prtesti #obs #p1 #p2 [, level(#)]`
- Using Stata as a *t*-test calculator, independent samples test (proportions): `prtesti #obs1 #p1 #obs2 #p2 [, level(#)]`

## Power calculation in Stata: two sample test

In Stata: Power analysis for a two-sample means test (independent samples).  $H_0$ : no difference. You select:

- Effect size (differences between groups)
- Equal variances or group-specific variances
- Sample size “allocation ratio”  $n_2/n_1$
- Significance level ( $\alpha$ )
- 2- or 1-sided test

As before, we can calculate power for given *ns* or determine the minimum *ns* needed to detect a given effect size at some level of power.

Note a “balanced design”  $n_2/n_1 = 1$  is more efficient (more power) than one with imbalanced sample sizes.

## Effect size

How should we express an *effect size* for a difference in means, in order to assess practical significance? Can use a Cohen's  $d$  type measure, expressing the estimated difference in means as a proportion of the pooled standard deviation:

$$d = \frac{\bar{x}_2 - \bar{x}_1}{s_p}$$

The `esize` command in Stata will calculate Cohen's  $d$  (and other effect size measures). It uses the pooled standard deviation calculation  $s_p$  (see earlier slides).

## Stata's esize

Stata's `esize` (effect size) command can calculate various types of effect sizes. For example, consider the  $t$ -test comparing 8th grade science achievement by gender:

```
. ttest achsci08,by(gender)
```

Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
male	226	58.34938	.6240243	9.381142	57.1197	59.57906
female	273	53.48392	.5340979	8.824746	52.43243	54.53541
combined	499	55.68751	.4203584	9.390095	54.86162	56.51341
diff		4.865461	.8166608		3.260928	6.469994

diff = mean(male) - mean(female) t = 5.9578  
Ho: diff = 0 degrees of freedom = 497

Ha: diff < 0  
Pr(T < t) = 1.0000

Ha: diff != 0  
Pr(|T| > |t|) = 0.0000

Ha: diff > 0  
Pr(T > t) = 0.0000

# Stata's esize

esize twosample can calculate Cohen's  $d$  (and other measures of effect size):

```
. esize twosample achsci08,by(gender)
```

Effect size based on mean comparison

Obs per group:  
male = 226  
female = 273

Effect Size	Estimate	[95% Conf. Interval]	
Cohen's $d$	.535793	.3561508	.714914
Hedges's $g$	.534984	.355613	.7138345

Note it uses the pooled standard deviation, not the “combined” std. dev. shown on the previous slide.