

Chapter 5: Inference for numerical data

OpenIntro Statistics, 3rd 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).

One-sample means with the t distribution

Friday the 13th

Between 1990 - 1992 researchers in the UK collected data on traffic flow, accidents, and hospital admissions on Friday 13th and the previous Friday, Friday 6th. Below is an excerpt from this data set on traffic flow. We can assume that traffic flow on given day at locations 1 and 2 are independent.

	type	date	6 th	13 th	diff	location
1	traffic	1990, July	139246	138548	698	loc 1
2	traffic	1990, July	134012	132908	1104	loc 2
3	traffic	1991, September	137055	136018	1037	loc 1
4	traffic	1991, September	133732	131843	1889	loc 2
5	traffic	1991, December	123552	121641	1911	loc 1
6	traffic	1991, December	121139	118723	2416	loc 2
7	traffic	1992, March	128293	125532	2761	loc 1
8	traffic	1992, March	124631	120249	4382	loc 2
9	traffic	1992, November	124609	122770	1839	loc 1
10	traffic	1992, November	117584	117263	321	loc 2

Friday the 13th

- We want to investigate if people's behavior is different on Friday 13th compared to Friday 6th.
- One approach is to compare the traffic flow on these two days.
- H_0 : Average traffic flow on Friday 6th and 13th are equal.
 H_A : Average traffic flow on Friday 6th and 13th are different.

Each case in the data set represents traffic flow recorded at the same location in the same month of the same year: one count from Friday 6th and the other Friday 13th. Are these two counts independent?

No

Hypotheses

What are the hypotheses for testing for a difference between the average traffic flow between Friday 6th and 13th?

(a) $H_0 : \mu_{6th} = \mu_{13th}$

$$H_A : \mu_{6th} \neq \mu_{13th}$$

(b) $H_0 : p_{6th} = p_{13th}$

$$H_A : p_{6th} \neq p_{13th}$$

(c) $H_0 : \mu_{diff} = 0$

$$H_A : \mu_{diff} \neq 0$$

(d) $H_0 : \bar{x}_{diff} = 0$

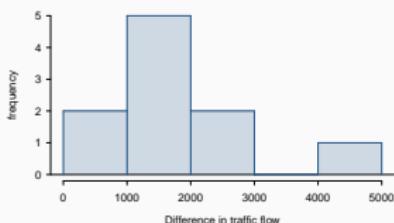
$$H_A : \bar{x}_{diff} \neq 0$$

Conditions

- *Independence:* We are told to assume that cases (rows) are independent.
- *Sample size / skew:*

Conditions

- *Independence:* We are told to assume that cases (rows) are independent.
- *Sample size / skew:*
- The sample distribution does not appear to be extremely skewed, but it's very difficult to assess with such a small sample size. We might want to think about whether we would expect the population distribution to be skewed or not – probably not, it should be equally likely to have days with lower than average traffic and higher than average traffic.
- We do not know σ and n is too small to assume s is a reliable estimate for σ .



So what do we do when the sample size is small?

Review: what purpose does a large sample serve?

As long as observations are independent, and the population distribution is not extremely skewed, a large sample would ensure that...

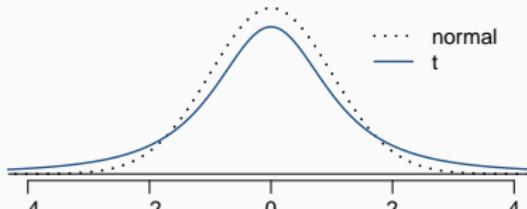
- the sampling distribution of the mean is nearly normal
- the estimate of the standard error, as $\frac{s}{\sqrt{n}}$, is reliable

The normality condition

- The CLT, which states that sampling distributions will be nearly normal, holds true for *any* sample size as long as the population distribution is nearly normal.
- While this is a helpful special case, it's inherently difficult to verify normality in small data sets.
- We should exercise caution when verifying the normality condition for small samples. It is important to not only examine the data but also think about where the data come from.
 - For example, ask: would I expect this distribution to be symmetric, and am I confident that outliers are rare?

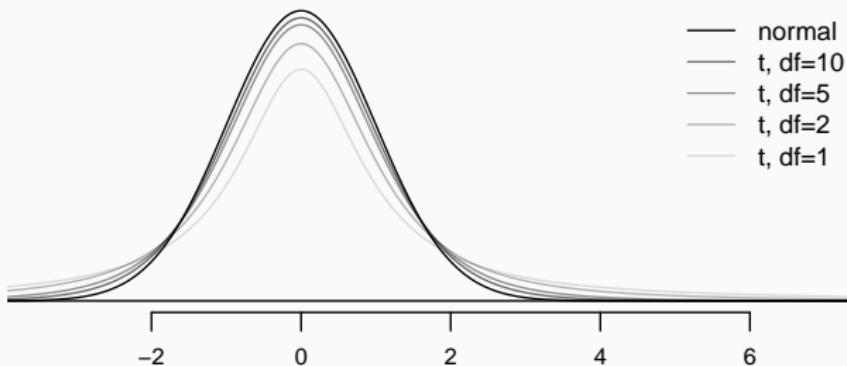
The t distribution

- When the population standard deviation is unknown (almost always), the uncertainty of the standard error estimate is addressed by using a new distribution: the t distribution.
- This distribution also has a bell shape, but its tails are *thicker* than the normal model's.
- Therefore observations are more likely to fall beyond two SDs from the mean than under the normal distribution.
- These extra thick tails are helpful for resolving our problem with a less reliable estimate the standard error (since n is small)



The t distribution (cont.)

- Always centered at zero, like the standard normal (z) distribution.
- Has a single parameter: *degrees of freedom (df)*.



What happens to shape of the t distribution as df increases?

Approaches normal.

Back to Friday the 13th

	type	date	6 th	13 th	diff	location
1	traffic	1990, July	139246	138548	698	loc 1
2	traffic	1990, July	134012	132908	1104	loc 2
3	traffic	1991, September	137055	136018	1037	loc 1
4	traffic	1991, September	133732	131843	1889	loc 2
5	traffic	1991, December	123552	121641	1911	loc 1
6	traffic	1991, December	121139	118723	2416	loc 2
7	traffic	1992, March	128293	125532	2761	loc 1
8	traffic	1992, March	124631	120249	4382	loc 2
9	traffic	1992, November	124609	122770	1839	loc 1
10	traffic	1992, November	117584	117263	321	loc 2



$$\bar{x}_{diff} = 1836$$

$$s_{diff} = 1176$$

Finding the test statistic

Test statistic for inference on a small sample mean

The test statistic for inference on a small sample ($n < 50$) mean is the T statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

in context...

$$\text{point estimate} = \bar{x}_{diff} = 1836$$

$$SE = \frac{s_{diff}}{\sqrt{n}} = \frac{1176}{\sqrt{10}} = 372$$

$$T = \frac{1836 - 0}{372} = 4.94$$

$$df = 10 - 1 = 9$$

Note: Null value is 0 because in the null hypothesis we set $\mu_{diff} = 0$.

Finding the p-value

- The p-value is, once again, calculated as the area tail area under the t distribution.
- Using R:

```
> 2 * pt(4.94, df = 9, lower.tail = FALSE)
```

```
[1] 0.0008022394
```

- Using a web app:
https://gallery.shinyapps.io/dist_calc/
- Or when these aren't available, we can use a t -table.

Finding the p-value

Locate the calculated T statistic on the appropriate df row, obtain the p-value from the corresponding column heading (one or two tail, depending on the alternative hypothesis).

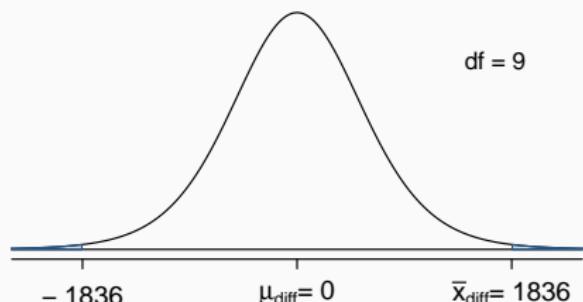
		one tail	0.100	0.050	0.025	0.010	0.005
		two tails	0.200	0.100	0.050	0.020	0.010
df	1		3.08	6.31	12.71	31.82	63.66
	2		1.89	2.92	4.30	6.96	9.92
	3		1.64	2.35	3.18	4.54	5.84
	⋮		⋮	⋮	⋮	⋮	⋮
	17		1.33	1.74	2.11	2.57	2.90
	18		1.33	1.73	2.10	2.55	2.88
	19		1.33	1.73	2.09	2.54	2.86
	20		1.33	1.72	2.09	2.53	2.85
	⋮		⋮	⋮	⋮	⋮	⋮
	400		1.28	1.65	1.97	2.34	2.59
	500		1.28	1.65	1.96	2.33	2.59
	∞		1.28	1.64	1.96	2.33	2.58

Finding the p-value (cont.)

one tail	0.100	0.050	0.025	0.010	0.005	
two tails	0.200	0.100	0.050	0.020	0.010	→
df	6	1.44	1.94	2.45	3.14	3.71
	7	1.41	1.89	2.36	3.00	3.50
	8	1.40	1.86	2.31	2.90	3.36
	9	1.38	1.83	2.26	2.82	3.25 →
	10	1.37	1.81	2.23	2.76	3.17

$$T = 4.94$$

What is the conclusion of the hypothesis test?



The data provide convincing evidence of a difference between traffic flow on Friday 6th and 13th.

What is the difference?

- We concluded that there is a difference in the traffic flow between Friday 6th and 13th.
- But it would be more interesting to find out what exactly this difference is.
- We can use a confidence interval to estimate this difference.

Confidence interval for a small sample mean

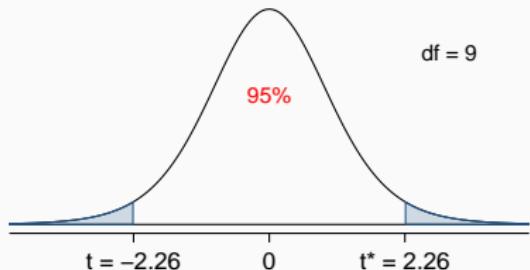
- Confidence intervals are always of the form

$$\text{point estimate} \pm ME$$

- ME is always calculated as the product of a critical value and SE.
- Since small sample means follow a t distribution (and not a z distribution), the critical value is a t^* (as opposed to a z^*).

$$\text{point estimate} \pm t^* \times SE$$

Finding the critical t (t^*)



$n = 10$, $df = 10 - 1 = 9$, t^* is at the intersection of row $df = 9$ and two tail probability 0.05.

one tail		0.100	0.050	0.025	0.010	0.005
two tails		0.200	0.100	0.050	0.020	0.010
df	6	1.44	1.94	2.45	3.14	3.71
	7	1.41	1.89	2.36	3.00	3.50
8	1.40	1.86	2.31	2.90	3.36	
9	1.38	1.83	2.26	2.82	3.25	
10	1.37	1.81	2.23	2.76	3.17	

Constructing a CI for a small sample mean

Which of the following is the correct calculation of a 95% confidence interval for the difference between the traffic flow between Friday 6th and 13th?

$$\bar{x}_{diff} = 1836 \quad s_{diff} = 1176 \quad n = 10 \quad SE = 372$$

- (a) $1836 \pm 1.96 \times 372$
- (b) $1836 \pm 2.26 \times 372 \rightarrow (995, 2677)$
- (c) $1836 \pm -2.26 \times 372$
- (d) $1836 \pm 2.26 \times 1176$

Interpreting the CI

Which of the following is the **best** interpretation for the confidence interval we just calculated?

$$\mu_{\text{diff}:6^{\text{th}}-13^{\text{th}}} = (995, 2677)$$

We are 95% confident that ...

- (a) the difference between the average number of cars on the road on Friday 6th and 13th is between 995 and 2,677.
- (b) on Friday 6th there are 995 to 2,677 fewer cars on the road than on the Friday 13th, on average.
- (c) on Friday 6th there are 995 fewer to 2,677 more cars on the road than on the Friday 13th, on average.
- (d) **on Friday 13th there are 995 to 2,677 fewer cars on the road than on the Friday 6th, on average.**

Synthesis

Does the conclusion from the hypothesis test agree with the findings of the confidence interval?

Yes, the hypothesis test found a significant difference, and the CI does not contain the null value of 0.

Do you think the findings of this study suggests that people believe Friday 13th is a day of bad luck?

No, this is an observational study. We have just observed a significant difference between the number of cars on the road on these two days. We have not tested for people's beliefs.

Recap: Inference using the t -distribution

- If σ is unknown, use the t -distribution with $SE = \frac{s}{\sqrt{n}}$.
- Conditions:
 - independence of observations (often verified by a random sample, and if sampling without replacement, $n < 10\%$ of population)
 - no extreme skew
- Hypothesis testing:

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}, \text{ where } df = n - 1$$

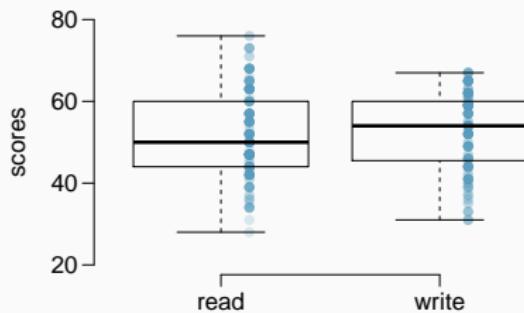
- Confidence interval:

$$\text{point estimate} \pm t_{df}^* \times SE$$

Note: The example we used was for paired means (difference between dependent groups). We took the difference between the observations and used only these differences (one sample) in our analysis, therefore the mechanics are

Paired data

200 observations were randomly sampled from the High School and Beyond survey. The same students took a reading and writing test and their scores are shown below. At a first glance, does there appear to be a difference between the average reading and writing test score?



The same students took a reading and writing test and their scores are shown below. Are the reading and writing scores of each student independent of each other?

	id	read	write
1	70	57	52
2	86	44	33
3	141	63	44
4	172	47	52
:	:	:	:
200	137	63	65

(a) Yes

(b) *No*

Analyzing paired data

- When two sets of observations have this special correspondence (not independent), they are said to be *paired*.
- To analyze paired data, it is often useful to look at the difference in outcomes of each pair of observations.

$$\text{diff} = \text{read} - \text{write}$$

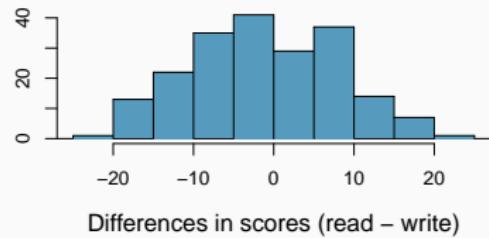
Analyzing paired data

- When two sets of observations have this special correspondence (not independent), they are said to be *paired*.
- To analyze paired data, it is often useful to look at the difference in outcomes of each pair of observations.

$$\text{diff} = \text{read} - \text{write}$$

- It is important that we always subtract using a consistent order.

	id	read	write	diff
1	70	57	52	5
2	86	44	33	11
3	141	63	44	19
4	172	47	52	-5
:	:	:	:	:
200	137	63	65	-2



Parameter and point estimate

- *Parameter of interest:* Average difference between the reading and writing scores of *all* high school students.

$$\mu_{diff}$$

- *Point estimate:* Average difference between the reading and writing scores of *sampled* high school students.

$$\bar{x}_{diff}$$

Setting the hypotheses

If in fact there was no difference between the scores on the reading and writing exams, what would you expect the average difference to be?

0

What are the hypotheses for testing if there is a difference between the average reading and writing scores?

H_0 : There is no difference between the average reading and writing score.

$$\mu_{diff} = 0$$

H_A : There is a difference between the average reading and writing score.

$$\mu_{diff} \neq 0$$

Nothing new here

- The analysis is no different than what we have done before.
- We have data from *one* sample: differences.
- We are testing to see if the average difference is different than 0.

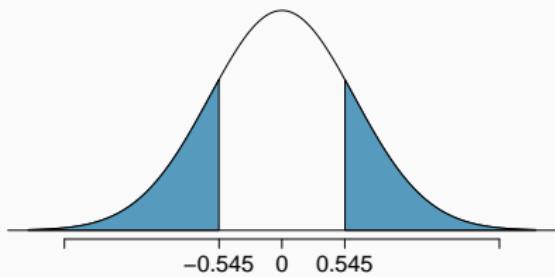
Checking assumptions & conditions

Which of the following is true?

- (a) *Since students are sampled randomly and are less than 10% of all high school students, we can assume that the difference between the reading and writing scores of one student in the sample is independent of another.*
- (b) The distribution of differences is bimodal, therefore we cannot continue with the hypothesis test.
- (c) In order for differences to be random we should have sampled with replacement.
- (d) Since students are sampled randomly and are less than 10% of all students, we can assume that the sampling distribution of the average difference will be nearly normal.

Calculating the test-statistic and the p-value

The observed average difference between the two scores is -0.545 points and the standard deviation of the difference is 8.887 points. Do these data provide convincing evidence of a difference between the average scores on the two exams? Use $\alpha = 0.05$.



$$T = \frac{-0.545 - 0}{\frac{8.887}{\sqrt{200}}} = \frac{-0.545}{0.628} = -0.87$$

$$df = 200 - 1 = 199$$

$$p\text{-value} = 0.1927 \times 2 = 0.3854$$

Since $p\text{-value} > 0.05$, fail to reject, the data do not provide convincing evidence of a difference between the average reading and writing scores.

Interpretation of p-value

Which of the following is the correct interpretation of the p-value?

- (a) Probability that the average scores on the reading and writing exams are equal.
- (b) Probability that the average scores on the reading and writing exams are different.
- (c) *Probability of obtaining a random sample of 200 students where the average difference between the reading and writing scores is at least 0.545 (in either direction), if in fact the true average difference between the scores is 0.*
- (d) Probability of incorrectly rejecting the null hypothesis if in fact the null hypothesis is true.

Suppose we were to construct a 95% confidence interval for the average difference between the reading and writing scores. Would you expect this interval to include 0?

- (a) yes
- (b) no
- (c) cannot tell from the information given

$$\begin{aligned}-0.545 \pm 1.97 \frac{8.887}{\sqrt{200}} &= -0.545 \pm 1.97 \times 0.628 \\ &= -0.545 \pm 1.24 \\ &= (-1.785, 0.695)\end{aligned}$$

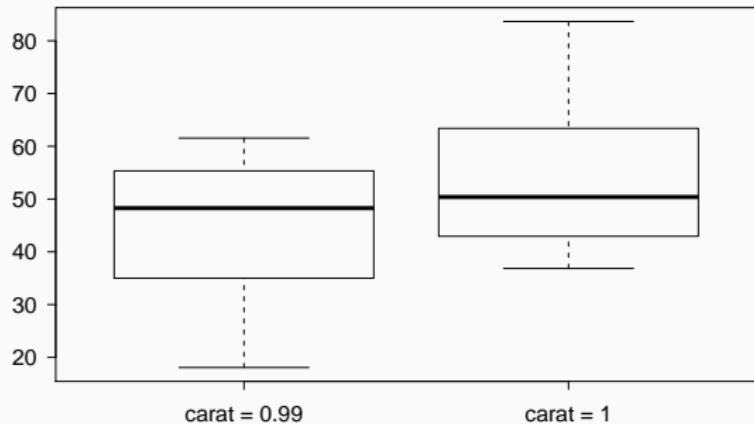
Difference of two means

Diamonds

- Weights of diamonds are measured in carats.
- 1 carat = 100 points, 0.99 carats = 99 points, etc.
- The difference between the size of a 0.99 carat diamond and a 1 carat diamond is undetectable to the naked human eye, but does the price of a 1 carat diamond tend to be higher than the price of a 0.99 diamond?
- We are going to test to see if there is a difference between the average prices of 0.99 and 1 carat diamonds.
- In order to be able to compare equivalent units, we divide the prices of 0.99 carat diamonds by 99 and 1 carat diamonds by 100, and compare the average point prices.



Data



	<i>0.99 carat</i>	<i>1 carat</i>
\bar{x}	44.50	53.43
s	13.32	12.22
n	23	30

These data are a random sample from the diamonds data set in ggplot2 R package.

Parameter and point estimate

- *Parameter of interest:* Average difference between the point prices of *all* 0.99 carat and 1 carat diamonds.

$$\mu_{pt99} - \mu_{pt100}$$

- *Point estimate:* Average difference between the point prices of *sampled* 0.99 carat and 1 carat diamonds.

$$\bar{x}_{pt99} - \bar{x}_{pt100}$$

Hypotheses

Which of the following is the correct set of hypotheses for testing if the average point price of 1 carat diamonds (μ_{pt100}) is higher than the average point price of 0.99 carat diamonds (μ_{pt99})?

- (a) $H_0 : \mu_{pt99} = \mu_{pt100}$
 $H_A : \mu_{pt99} \neq \mu_{pt100}$
- (b) $H_0 : \mu_{pt99} = \mu_{pt100}$
 $H_A : \mu_{pt99} > \mu_{pt100}$
- (c) $H_0 : \mu_{pt99} = \mu_{pt100}$
 $H_A : \mu_{pt99} < \mu_{pt100}$
- (d) $H_0 : \bar{x}_{pt99} = \bar{x}_{pt100}$
 $H_A : \bar{x}_{pt99} < \bar{x}_{pt100}$

Conditions

Which of the following does not need to be satisfied in order to conduct this hypothesis test using theoretical methods?

- (a) Point price of one 0.99 carat diamond in the sample should be independent of another, and the point price of one 1 carat diamond should independent of another as well.
- (b) Point prices of 0.99 carat and 1 carat diamonds in the sample should be independent.
- (c) Distributions of point prices of 0.99 and 1 carat diamonds should not be extremely skewed.
- (d) *Both sample sizes should be at least 30.*

Test statistic

Test statistic for inference on the difference of two small sample means

The test statistic for inference on the difference of two means where σ_1 and σ_2 are unknown is the T statistic.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

where

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad \text{and} \quad df = \min(n_1 - 1, n_2 - 1)$$

Note: The calculation of the df is actually much more complicated. For simplicity we'll use the above formula to estimate the true df when conducting the analysis by hand.

Test statistic (cont.)

	0.99 carat pt99	1 carat pt100
\bar{x}	44.50	53.43
s	13.32	12.22
n	23	30

in context...

$$\begin{aligned} T &= \frac{\text{point estimate} - \text{null value}}{SE} \\ &= \frac{(44.50 - 53.43) - 0}{\sqrt{\frac{13.32^2}{23} + \frac{12.22^2}{30}}} \\ &= \frac{-8.93}{3.56} \\ &= -2.508 \end{aligned}$$

Test statistic (cont.)

Which of the following is the correct df for this hypothesis test?

- (a) 22 $\rightarrow df = \min(n_{pt99} - 1, n_{pt100} - 1)$
- (b) 23 $= \min(23 - 1, 30 - 1)$
- (c) 30 $= \min(22, 29) = 22$
- (d) 29
- (e) 52

p-value

Which of the following is the correct p-value for this hypothesis test?

$$T = -2.508 \quad df = 22$$

- (a) between 0.005 and 0.01
- (b) *between 0.01 and 0.025*
- (c) between 0.02 and 0.05
- (d) between 0.01 and 0.02

		one tail	0.100	0.050	0.025	0.010
		two tails	0.200	0.100	0.050	0.020
df	21	1.32	1.72	2.08	2.52	
	22	1.32	1.72	2.07	2.51	
	23	1.32	1.71	2.07	2.50	
	24	1.32	1.71	2.06	2.49	
	25	1.32	1.71	2.06	2.49	

Synthesis

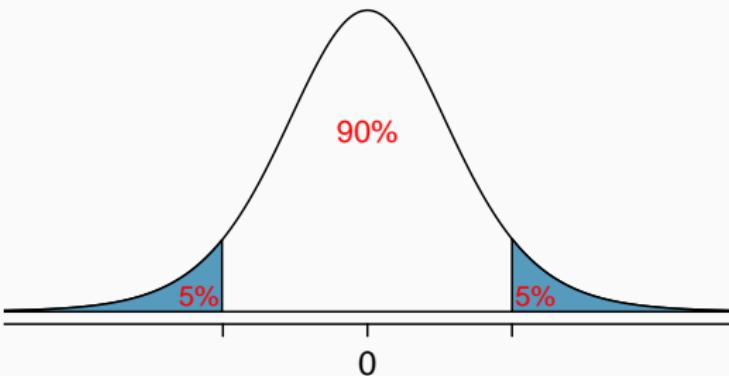
What is the conclusion of the hypothesis test? How (if at all) would this conclusion change your behavior if you went diamond shopping?

- *p-value is small so reject H_0 . The data provide convincing evidence to suggest that the point price of 0.99 carat diamonds is lower than the point price of 1 carat diamonds.*
- *Maybe buy a 0.99 carat diamond? It looks like a 1 carat, but is significantly cheaper.*

Equivalent confidence level

What is the equivalent confidence level for a one-sided hypothesis test at $\alpha = 0.05$?

- (a) 90%
- (b) 92.5%
- (c) 95%
- (d) 97.5%



Critical value

What is the appropriate t^* for a confidence interval for the average difference between the point prices of 0.99 and 1 carat diamonds?

- (a) 1.32
- (b) 1.72
- (c) 2.07
- (d) 2.82

one tail		0.100	0.050	0.025	0.010	0.005
two tails		0.200	0.100	0.050	0.020	0.010
df	21	1.32	1.72	2.08	2.52	2.83
	22	1.32	1.72	2.07	2.51	2.82
	23	1.32	1.71	2.07	2.50	2.81
	24	1.32	1.71	2.06	2.49	2.80
	25	1.32	1.71	2.06	2.49	2.79

Confidence interval

Calculate the interval, and interpret it in context.

Confidence interval

Calculate the interval, and interpret it in context.

$$\text{point estimate} \pm ME$$

$$\begin{aligned} (\bar{x}_{pt99} - \bar{x}_{pt1}) \pm t_{df}^{\star} \times SE &= (44.50 - 53.43) \pm 1.72 \times 3.56 \\ &= -8.93 \pm 6.12 \\ &= (-15.05, -2.81) \end{aligned}$$

We are 90% confident that the average point price of a 0.99 carat diamond is \$15.05 to \$2.81 lower than the average point price of a 1 carat diamond.

Recap: Inference using difference of two small sample means

- If σ_1 or σ_2 is unknown, difference between the sample means follow a t -distribution with $SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_1}}$.
- Conditions:
 - independence within groups (often verified by a random sample, and if sampling without replacement, $n < 10\%$ of population) and between groups
 - no extreme skew in either group
- Hypothesis testing:
$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}, \text{ where } df = \min(n_1 - 1, n_2 - 1)$$
- Confidence interval:

$$\text{point estimate} \pm t_{df}^* \times SE$$

Computing the power for a 2-sample test

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

- Type 1 error is rejecting H_0 when you shouldn't have, and the probability of doing so is α (significance level)
- Type 2 error is failing to reject H_0 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 $1 - \beta$
- 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 H_0).
- 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 (BP), hypotheses

Suppose a pharmaceutical company has developed a new drug for lowering blood pressure, and they are preparing a clinical trial to test the drug's effectiveness. They recruit people who are taking a particular standard blood pressure medication, and half of the subjects are given the new drug (treatment) and the other half continue to take their current medication through generic-looking pills to ensure blinding (control). What are the hypotheses for a two-sided hypothesis test in this context?

$$H_0 : \mu_{treatment} - \mu_{control} = 0$$

$$H_A : \mu_{treatment} - \mu_{control} \neq 0$$

Example - BP, standard error

Suppose researchers would like to run the clinical trial on patients with systolic blood pressures between 140 and 180 mmHg. Suppose previously published studies suggest that the standard deviation of the patients' blood pressures will be about 12 mmHg and the distribution of patient blood pressures will be approximately symmetric. If we had 100 patients per group, what would be the approximate standard error for difference in sample means of the treatment and control groups?

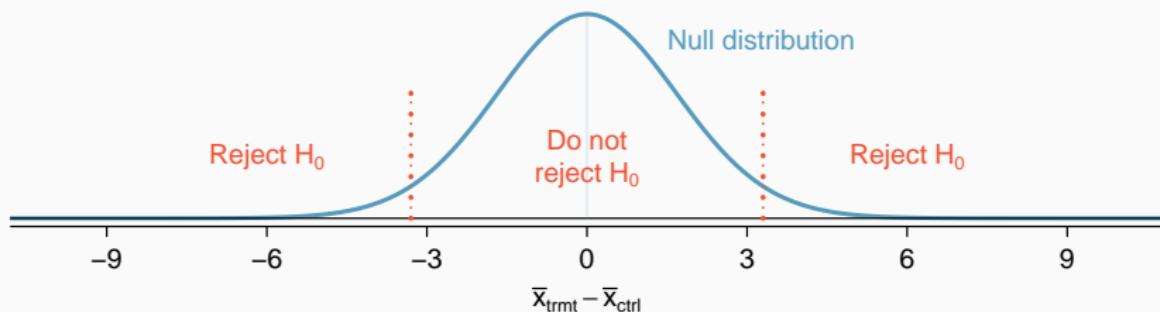
$$SE = \sqrt{\frac{12^2}{100} + \frac{12^2}{100}} = 1.70$$

Example - BP, minimum effect size required to reject H_0

For what values of the difference between the observed averages of blood pressure in treatment and control groups (effect size) would we reject the null hypothesis at the 5% significance level?

Example - BP, minimum effect size required to reject H_0

For what values of the difference between the observed averages of blood pressure in treatment and control groups (effect size) would we reject the null hypothesis at the 5% significance level?



The difference should be at least

$$1.96 * 1.70 = 3.332$$

or at most

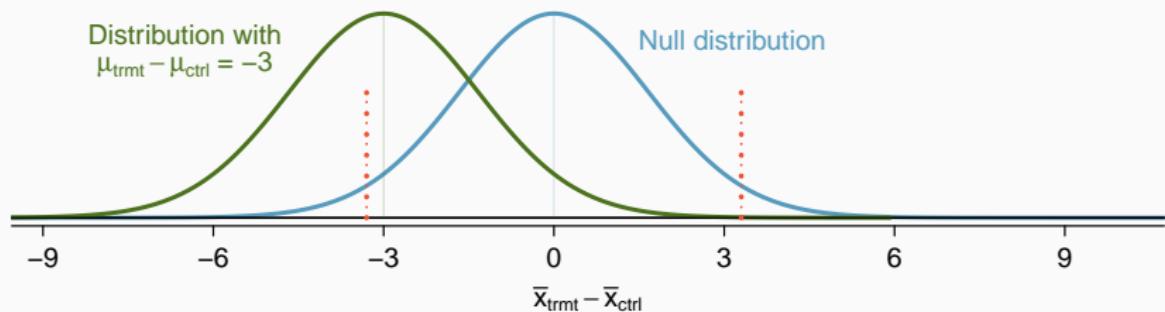
$$-1.96 * 1.70 = 3.332.$$

Example - BP, power

Suppose that the company researchers care about finding any effect on blood pressure that is 3 mmHg or larger vs the standard medication. What is the power of the test that can detect this effect?

Example - BP, power

Suppose that the company researchers care about finding any effect on blood pressure that is 3 mmHg or larger vs the standard medication. What is the power of the test that can detect this effect?



$$Z = \frac{-3.332 - (-3)}{1.70} = -0.20$$

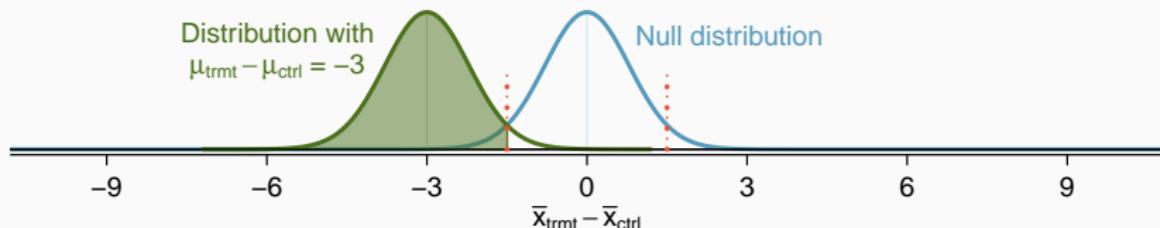
$$P(Z < -0.20) = 0.4207$$

Example - BP, required sample size for 80% power

What sample size will lead to a power of 80% for this test?

Example - BP, required sample size for 80% power

What sample size will lead to a power of 80% for this test?



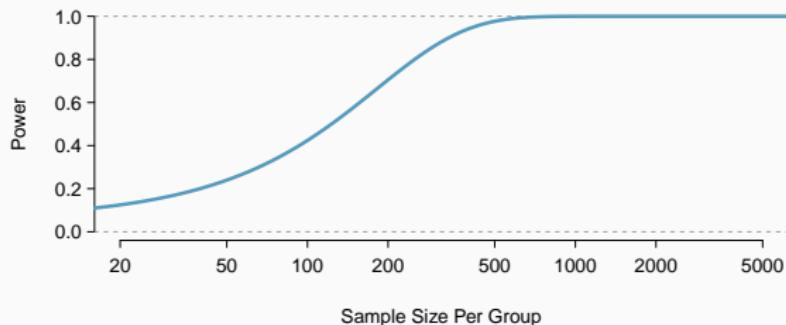
$$SE = \frac{3}{2.8} = 1.07142$$

$$1.07142 = \sqrt{\frac{12^2}{n} + \frac{12^2}{n}}$$

$$n = 250.88 \rightarrow n \geq 251$$

Recap

- Calculate required sample size for a desired level of power
- Calculate power for a range of sample sizes, then choose the sample size that yields the target power (usually 80% or 90%)



Achieving desired power

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.

Comparing means with ANOVA



- The Wolf River in Tennessee flows past an abandoned site once used by the pesticide industry for dumping wastes, including chlordane (pesticide), aldrin, and dieldrin (both insecticides).
- These highly toxic organic compounds can cause various cancers and birth defects.
- The standard methods to test whether these substances are present in a river is to take samples at six-tenths depth.
- But since these compounds are denser than water and their molecules tend to stick to particles of sediment, they are more likely to be found in higher concentrations near the bottom than near

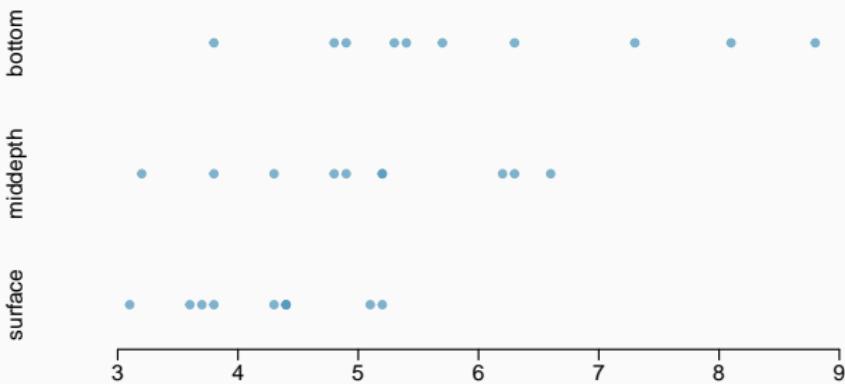
Data

Aldrin concentration (nanograms per liter) at three levels of depth.

	aldrin	depth
1	3.80	bottom
2	4.80	bottom
...		
10	8.80	bottom
11	3.20	middepth
12	3.80	middepth
...		
20	6.60	middepth
21	3.10	surface
22	3.60	surface
...		
30	5.20	surface

Exploratory analysis

Aldrin concentration (nanograms per liter) at three levels of depth.



	n	mean	sd
bottom	10	6.04	1.58
middepth	10	5.05	1.10
surface	10	4.20	0.66
overall	30	5.10	1.37

Research question

Is there a difference between the mean aldrin concentrations among the three levels?

- To compare means of 2 groups we use a Z or a T statistic.
- To compare means of 3+ groups we use a new test called **ANOVA** and a new statistic called **F**.

ANOVA

ANOVA is used to assess whether the mean of the outcome variable is different for different levels of a categorical variable.

H_0 : The mean outcome is the same across all categories,

$$\mu_1 = \mu_2 = \cdots = \mu_k,$$

where μ_i represents the mean of the outcome for observations in category i .

H_A : At least one mean is different than others.

Conditions

1. The observations should be independent within and between groups
 - If the data are a simple random sample from less than 10% of the population, this condition is satisfied.
 - Carefully consider whether the data may be independent (e.g. no pairing).
 - Always important, but sometimes difficult to check.
2. The observations within each group should be nearly normal.
 - Especially important when the sample sizes are small.

How do we check for normality?

3. The variability across the groups should be about equal.
 - Especially important when the sample sizes differ between groups.

How can we check this condition?

z/t test vs. ANOVA - Purpose

z/t test

Compare means from *two* groups to see whether they are so far apart that the observed difference cannot reasonably be attributed to sampling variability.

$$H_0 : \mu_1 = \mu_2$$

ANOVA

Compare the means from *two or more* groups to see whether they are so far apart that the observed differences cannot all reasonably be attributed to sampling variability.

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

z/t test vs. ANOVA - Method

z/t test

Compute a test statistic (a ratio).

$$z/t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{SE(\bar{x}_1 - \bar{x}_2)}$$

ANOVA

Compute a test statistic (a ratio).

$$F = \frac{\text{variability bet. groups}}{\text{variability w/in groups}}$$

- Large test statistics lead to small p-values.
- If the p-value is small enough H_0 is rejected, we conclude that the population means are not equal.

z/t test vs. ANOVA

- With only two groups t-test and ANOVA are equivalent, but only if we use a pooled standard variance in the denominator of the test statistic.
- With more than two groups, ANOVA compares the sample means to an overall *grand mean*.

Hypotheses

What are the correct hypotheses for testing for a difference between the mean aldrin concentrations among the three levels?

(a) $H_0 : \mu_B = \mu_M = \mu_S$

$$H_A : \mu_B \neq \mu_M \neq \mu_S$$

(b) $H_0 : \mu_B \neq \mu_M \neq \mu_S$

$$H_A : \mu_B = \mu_M = \mu_S$$

(c) $H_0 : \mu_B = \mu_M = \mu_S$

H_A : At least one mean is different.

(d) $H_0 : \mu_B = \mu_M = \mu_S = 0$

H_A : At least one mean is different.

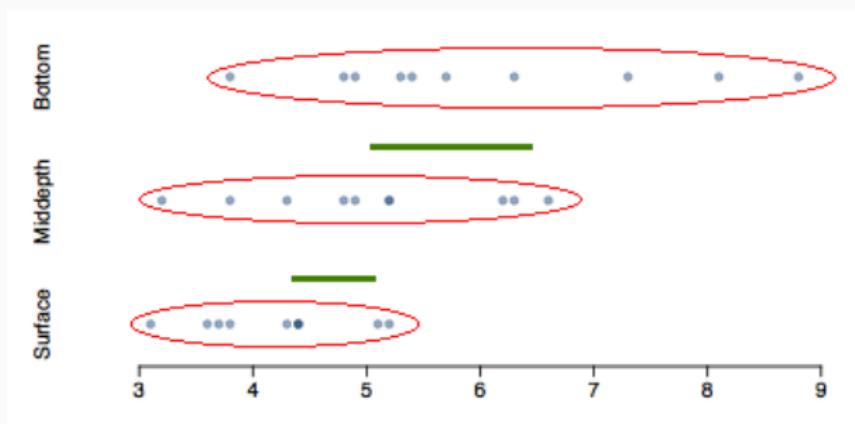
(e) $H_0 : \mu_B = \mu_M = \mu_S$

$$H_A : \mu_B > \mu_M > \mu_S$$

Test statistic

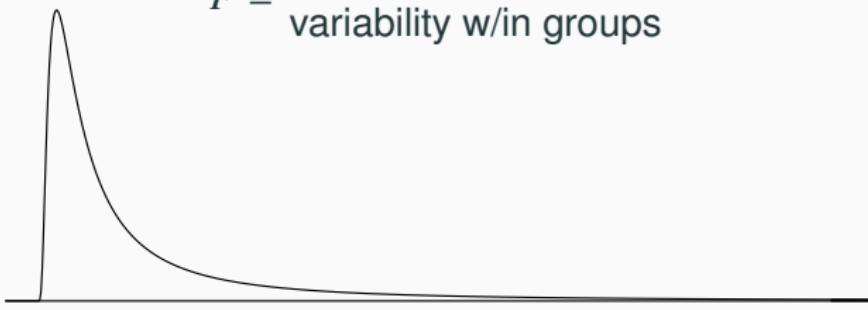
Does there appear to be a lot of variability within groups? How about between groups?

$$F = \frac{\text{variability bet. groups}}{\text{variability w/in groups}}$$



F distribution and p-value

$$F = \frac{\text{variability bet. groups}}{\text{variability w/in groups}}$$



- In order to be able to reject H_0 , we need a small p-value, which requires a large F statistic.
- In order to obtain a large F statistic, variability between sample means needs to be greater than variability within sample means.

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

Degrees of freedom associated with ANOVA

- groups: $df_G = k - 1$, where k is the number of groups
- total: $df_T = n - 1$, where n is the total sample size
- error: $df_E = df_T - df_G$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

Degrees of freedom associated with ANOVA

- groups: $df_G = k - 1$, where k is the number of groups
- total: $df_T = n - 1$, where n is the total sample size
- error: $df_E = df_T - df_G$

- $df_G = k - 1 = 3 - 1 = 2$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

Degrees of freedom associated with ANOVA

- groups: $df_G = k - 1$, where k is the number of groups
- total: $df_T = n - 1$, where n is the total sample size
- error: $df_E = df_T - df_G$

- $df_G = k - 1 = 3 - 1 = 2$
- $df_T = n - 1 = 30 - 1 = 29$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

Degrees of freedom associated with ANOVA

- groups: $df_G = k - 1$, where k is the number of groups
- total: $df_T = n - 1$, where n is the total sample size
- error: $df_E = df_T - df_G$

- $df_G = k - 1 = 3 - 1 = 2$
- $df_T = n - 1 = 30 - 1 = 29$
- $df_E = 29 - 2 = 27$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

Sum of squares between groups, SSG

Measures the variability between groups

$$SSG = \sum_{i=1}^k n_i(\bar{x}_i - \bar{x})^2$$

where n_i is each group size, \bar{x}_i is the average for each group, \bar{x} is the overall (grand) mean.

	n	mean
bottom	10	6.04
middepth	10	5.05
surface	10	4.2
overall	30	5.1

$$\begin{aligned} SSG &= (10 \times (6.04 - 5.1)^2) \\ &\quad + (10 \times (5.05 - 5.1)^2) \\ &\quad + (10 \times (4.2 - 5.1)^2) \\ &= 16.96 \end{aligned}$$

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group) depth	2	16.96	8.48	6.13	0.0063
(Error) Residuals	27	37.33	1.38		
Total	29	54.29			

Sum of squares total, SST

Measures the variability between groups

$$SST = \sum_{i=1}^n (x_i - \bar{x})$$

where x_i represent each observation in the dataset.

$$\begin{aligned}
 SST &= (3.8 - 5.1)^2 + (4.8 - 5.1)^2 + (4.9 - 5.1)^2 + \dots + (5.2 - 5.1)^2 \\
 &= (-1.3)^2 + (-0.3)^2 + (-0.2)^2 + \dots + (0.1)^2 \\
 &= 1.69 + 0.09 + 0.04 + \dots + 0.01 \\
 &= 54.29
 \end{aligned}$$

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group) depth	2	16.96	8.48	6.13	0.0063
(Error) Residuals	27	37.33	1.38		
Total	29	54.29			

Sum of squares error, SSE

Measures the variability within groups:

$$SSE = SST - SSG$$

$$SSE = 54.29 - 16.96 = 37.33$$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.13	0.0063
(Error)	Residuals	27	37.33	1.38		
	Total	29	54.29			

Mean square error

Mean square error is calculated as sum of squares divided by the degrees of freedom.

$$MSG = 16.96/2 = 8.48$$

$$MSE = 37.33/27 = 1.38$$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.14	0.0063
(Error)	Residuals	27	37.33	1.38		
	Total	29	54.29			

Test statistic, F value

As we discussed before, the F statistic is the ratio of the between group and within group variability.

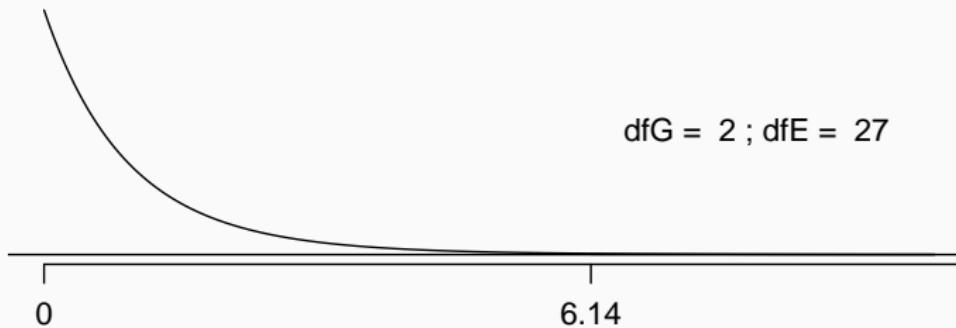
$$F = \frac{MSG}{MSE}$$

$$F = \frac{8.48}{1.38} = 6.14$$

		Df	Sum Sq	Mean Sq	F value	Pr(>F)
(Group)	depth	2	16.96	8.48	6.14	0.0063
Error	Residuals	27	37.33	1.38		
	Total	29	54.29			

p-value

p-value is the probability of at least as large a ratio between the “between group” and “within group” variability, if in fact the means of all groups are equal. It’s calculated as the area under the F curve, with degrees of freedom df_G and df_E , above the observed F statistic.



Conclusion - in context

What is the conclusion of the hypothesis test?

The data provide convincing evidence that the average aldrin concentration

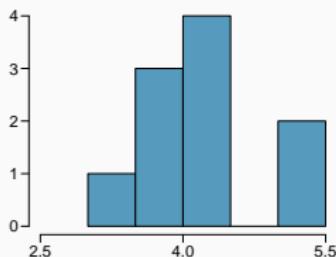
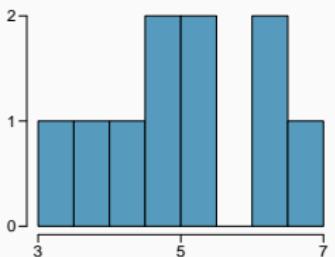
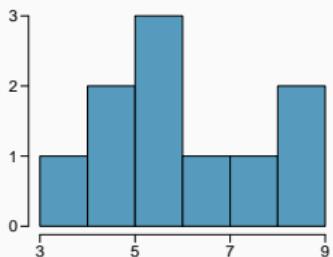
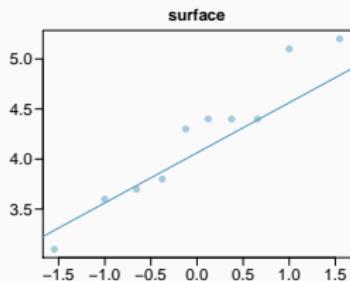
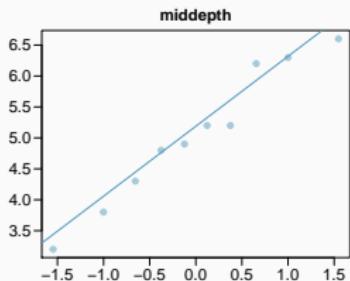
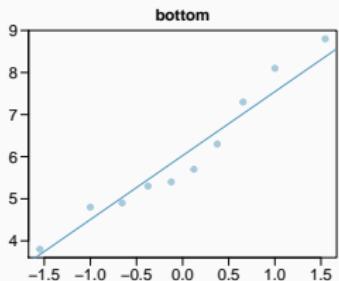
- (a) is different for all groups.
- (b) on the surface is lower than the other levels.
- (c) *is different for at least one group.*
- (d) is the same for all groups.

Conclusion

- If p-value is small (less than α), reject H_0 . The data provide convincing evidence that at least one mean is different from (but we can't tell which one).
- If p-value is large, fail to reject H_0 . The data do not provide convincing evidence that at least one pair of means are different from each other, the observed differences in sample means are attributable to sampling variability (or chance).

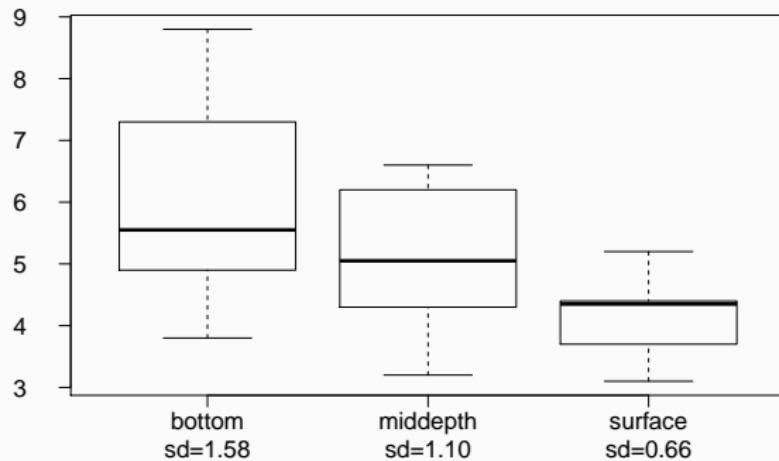
(2) approximately normal

Does this condition appear to be satisfied?



(3) constant variance

Does this condition appear to be satisfied?



Which means differ?

- Earlier we concluded that at least one pair of means differ.
The natural question that follows is “which ones?”
- We can do two sample t tests for differences in each possible pair of groups.

Can you see any pitfalls with this approach?

- When we run too many tests, the Type 1 Error rate increases.
- This issue is resolved by using a modified significance level.

Multiple comparisons

- The scenario of testing many pairs of groups is called *multiple comparisons*.
- The *Bonferroni correction* suggests that a more *stringent* significance level is more appropriate for these tests:

$$\alpha^* = \alpha/K$$

where K is the number of comparisons being considered.

- If there are k groups, then usually all possible pairs are compared and $K = \frac{k(k-1)}{2}$.

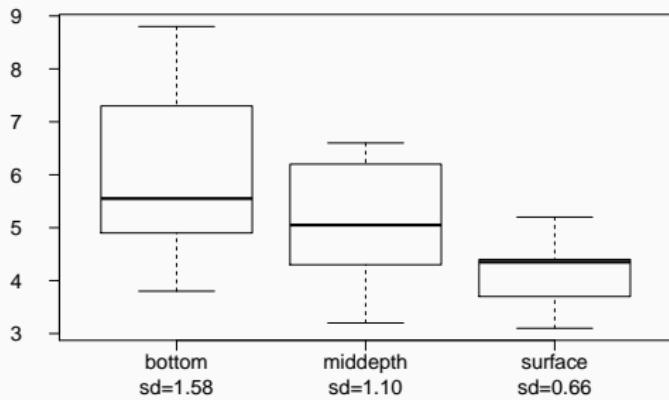
Determining the modified α

In the aldrin data set depth has 3 levels: bottom, mid-depth, and surface. If $\alpha = 0.05$, what should be the modified significance level for two sample t tests for determining which pairs of groups have significantly different means?

- (a) $\alpha^* = 0.05$
- (b) $\alpha^* = 0.05/2 = 0.025$
- (c) $\alpha^* = 0.05/3 = 0.0167$
- (d) $\alpha^* = 0.05/6 = 0.0083$

Which means differ?

Based on the box plots below, which means would you expect to be significantly different?



- (a) bottom & surface
- (b) bottom & mid-depth
- (c) mid-depth & surface
- (d) bottom & mid-depth;
mid-depth & surface
- (e) bottom & mid-depth;
bottom & surface;
mid-depth & surface

Which means differ? (cont.)

If the ANOVA assumption of equal variability across groups is satisfied, we can use the data from all groups to estimate variability:

- Estimate any within-group standard deviation with \sqrt{MSE} , which is s_{pooled}
- Use the error degrees of freedom, $n - k$, for t -distributions

Difference in two means: after ANOVA

$$SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{MSE}{n_1} + \frac{MSE}{n_2}}$$

Is there a difference between the average aldrin concentration at the bottom and at mid depth?

	n	mean	sd	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bottom	10	6.04	1.58					
middepth	10	5.05	1.10	depth	16.96	8.48	6.13	0.0063
surface	10	4.2	0.66	Residuals	27	37.33	1.38	
overall	30	5.1	1.37	Total	29	54.29		

$$T_{df_E} = \frac{(\bar{x}_{bottom} - \bar{x}_{middepth})}{\sqrt{\frac{MSE}{n_{bottom}} + \frac{MSE}{n_{middepth}}}}$$

$$T_{27} = \frac{(6.04 - 5.05)}{\sqrt{\frac{1.38}{10} + \frac{1.38}{10}}} = \frac{0.99}{0.53} = 1.87$$

$$0.05 < p-value < 0.10 \quad (\text{two-sided})$$

$$\alpha^* = 0.05/3 = 0.0167$$

Fail to reject H_0 , the data do not provide convincing evidence of a difference between the average aldrin concentrations at bottom and mid depth

Pairwise comparisons

Is there a difference between the average aldrin concentration at the bottom and at surface?

$$T_{df_E} = \frac{(\bar{x}_{bottom} - \bar{x}_{surface})}{\sqrt{\frac{MSE}{n_{bottom}} + \frac{MSE}{n_{surface}}}}$$

$$T_{27} = \frac{(6.04 - 4.02)}{\sqrt{\frac{1.38}{10} + \frac{1.38}{10}}} = \frac{2.02}{0.53} = 3.81$$

$$p-value < 0.01 \quad (two-sided)$$

$$\alpha^* = 0.05/3 = 0.0167$$

Reject H_0 , the data provide convincing evidence of a difference between the average aldrin concentrations at bottom and surface.