

Sta 111 - Summer II 2017
Probability and Statistical Inference
12. ANOVA

Lu Wang

Duke University, Department of Statistical Science

July 21, 2017

Outline

1. Comparing means of many different groups with ANOVA (analysis of variance)
2. ANOVA compares between group variation to within group variation
 1. ANOVA table
3. To identify which means are different, use t -tests and the Bonferroni correction
 1. Use a modified significance level in multiple comparisons
 2. The pooled standard deviation estimate from ANOVA

Aldrin in the Wolf River



- ▶ 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).
- ▶ The goal is to test whether these substances are present in a river by taking samples at different depths.
- ▶ 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 mid-depth.

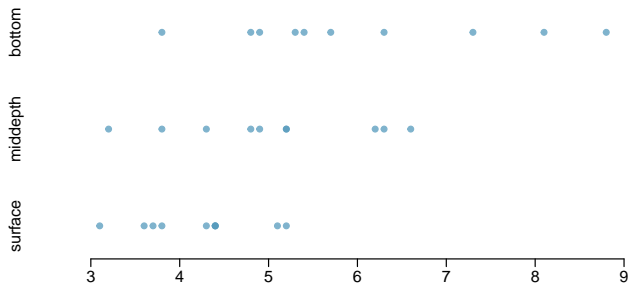
Data

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 is used to assess whether the mean of the outcome variable is different for different categories.

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

$$\mu_1 = \mu_2 = \dots = \mu_k,$$

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

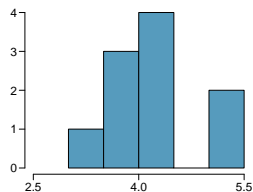
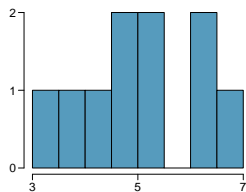
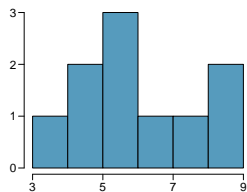
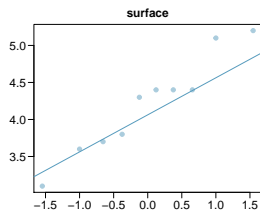
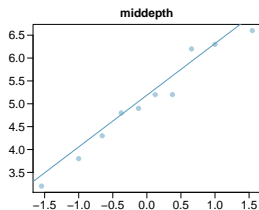
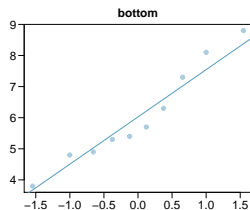
H_A : At least one mean is different than others.

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).
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?

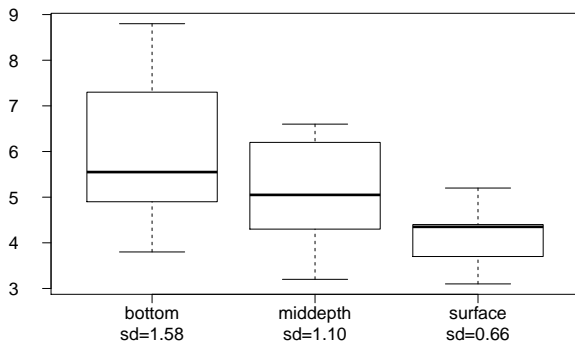
Check conditions

Does the "approximately normal" condition appear to be satisfied?



Check conditions

Does the "constant variance" condition appear to be satisfied?



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

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

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).

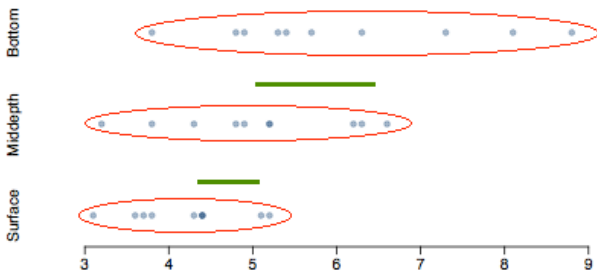
$$F = \frac{\text{variability bet. groups}}{\text{variability w/in groups}}$$

- ▶ 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.

Test statistic

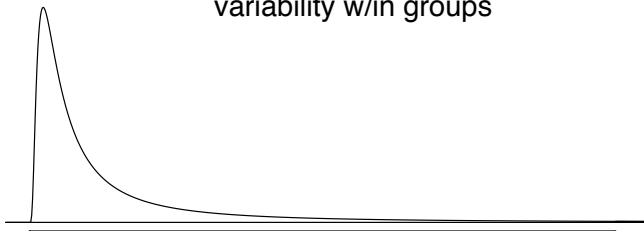
Does there appear to be a lot of variability within groups? How about between groups?

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



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: $df_E = df_T - df_G$

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

		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	$SSG = (10 \times (6.04 - 5.1)^2)$ $+ (10 \times (5.05 - 5.1)^2)$ $+ (10 \times (4.2 - 5.1)^2)$ $= 16.96$
middepth	10	5.05	
surface	10	4.2	
overall	30	5.1	

		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 overall variability

$$SST = \sum_{i=1}^n (x_i - \bar{x})^2$$

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

Mean square is calculated as sum of squares divided by the associated 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.

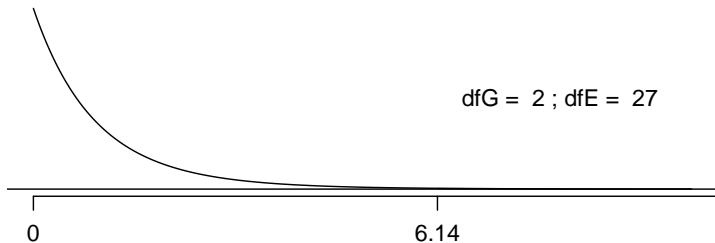
$$F = \frac{MSG}{MSE}$$

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

		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			

F test and p-value

If H_0 is true (the means of all groups are equal) and the model assumptions are satisfied, the F statistic follows an F distribution with parameters $df_1 = df_G$ and $df_2 = df_E$. p-value is calculated as the area under the F curve and above the observed F statistic (upper tail).



What is the conclusion of the hypothesis test?

The data provide significant 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 significant evidence that at least one mean is different (but we can't tell which one).

- ▶ If p-value is large, fail to reject H_0 . The data do not provide significant evidence that at least one mean is different. The observed differences in sample means are likely due to sampling variability (or chance).

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 a two-sample t test for each 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 in determining which pairs of groups have significantly different means?

(a) $\alpha^* = 0.05$

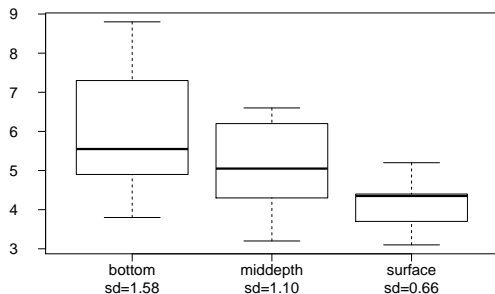
(b) $\alpha^* = 0.05/2 = 0.025$

(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 make the t -distribution approach slightly more precise by using a *pooled standard deviation*:

- ▶ The pooled standard deviation is a way to use data from all the groups to better estimate the standard deviation for each group.
- ▶ By pooling all the data, we can use a larger degree of freedom for the t -distribution.
- ▶ Both of these changes may permit a more accurate model of the sampling distribution of $\bar{x}_1 - \bar{x}_2$ if the standard deviations of the groups are equal.

Pooled standard deviation estimate from ANOVA

- ▶ The standard deviation of each group is estimated as
 $S_{pooled} = \sqrt{MSE}$
- ▶ Use the error degrees of freedom, $n - k$, for t -distributions
- ▶ The standard error of test statistic

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

Pairwise comparison: bottom vs mid depth

Suppose we know that $\sigma_1 = \sigma_2 = \sigma_3$. Is there a difference between the average aldrin concentration at the bottom and at mid depth?

bottom	10	6.04	1.58					
middepth	10	5.05	1.10					
surface	10	4.2	0.66					
overall	30	5.1	1.37					
				Df	Sum Sq	Mean Sq	F value	Pr(>F)
				depth	2	8.48	6.13	0.0063
				Residuals	27	1.38		
				Total	29	54.29		

$$T_{df_E} = \frac{(\bar{X}_{bottom} - \bar{X}_{middepth})}{\sqrt{\frac{MSE}{n_{bottom}} + \frac{MSE}{n_{middepth}}}}$$
$$T_{27} = \frac{(6.04 - 5.05)}{\sqrt{\frac{1.38}{10} + \frac{1.38}{10}}} = \frac{0.99}{0.53} = 1.87$$

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

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

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: bottom vs surface

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

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

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

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

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

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