

2.3: Tukey Test for Pairwise Mean Comparisons

If (and only if) we reject the null hypothesis, we then conclude at least one group is different from one other (importantly, we do *not* conclude that all the groups are different).

If it is the case that we reject the null, then we will want to know *which* group or groups are different. In our example we are not satisfied knowing at least one treatment level is different, we want to know where the difference is and the nature of the difference. To answer this question, we can follow up the ANOVA with a mean comparison procedure to find out which means differ from each other and which ones don't.

You might think we could not bother with the ANOVA and proceed with a series of t-tests to compare the groups. While that is intuitively simple, it creates inflation of the type I error. How does this inflation of type I error happen? For a single test,

$$\alpha = 1 - (.95) \quad (2.3.1)$$

The probability of committing a type I error (by random chance) for two simultaneous tests follows from the Multiplication Rule for independent events in probability. Recall that, for two independent events A and B the probability of A and B both occurring is $P(A \text{ and } B) = P(A) * P(B)$. So for two tests, we have

$$\alpha = 1 - ((.95) * (.95)) = 0.0975 \quad (2.3.2)$$

which is now larger than the α that we originally set. For our example, we have 6 comparisons, so $\alpha = 1 - (.95^6) = 0.2649$ which is a much larger (inflated) probability of committing a type I error than we originally set.

The multiple comparison procedures compensate for the type I error inflation (although each does so in a slightly different way).

There are several comparison procedures that can be employed, but we will start with the one most commonly used, the Tukey procedure. In the Tukey procedure, we compute a "yardstick" value based on the MS_{Error} and the number of means being compared. If any two means differ by more than the Tukey w value, then they are significantly different.

Step 1: Compute Tukey's w value

$$w = q_{\alpha(p, df_{Error})} \cdot s_{\bar{Y}} \quad (2.3.3)$$

where q_{α} is obtained from a table of Tukey q values

p = the number of treatment levels

$s_{\bar{Y}}$ = standard error of a treatment mean = $\sqrt{MS_{Error}/r}$

r = number of replications

Show Tukey q Values Table

df for Error Term	α	p = Number of Treatments								
		2	3	4	5	6	7	8	9	10
5	0.05	3.64	4.6	5.22	5.67	6.03	6.33	6.58	6.80	6.99
	0.01	5.70	6.98	7.80	8.42	8.91	9.32	9.67	9.97	10.24
6	0.05	3.46	4.34	4.90	5.30	5.63	5.90	6.12	6.32	6.49
	0.01	5.24	6.33	7.03	7.56	7.97	8.32	8.61	8.87	9.10
7	0.05	3.34	4.16	4.68	5.06	5.36	5.61	5.82	6.00	6.16
	0.01	4.95	5.92	6.54	7.01	7.37	7.68	7.94	8.17	8.37
8	0.05	3.26	4.04	4.53	4.89	5.17	5.40	5.60	5.77	5.92
	0.01	4.75	5.64	6.20	6.62	6.96	7.24	7.47	7.68	7.86
9	0.05	3.20	3.95	4.41	4.76	5.02	5.24	5.43	5.59	5.74
	0.01	4.60	5.43	5.96	6.35	6.66	6.91	7.13	7.33	7.49

df for Error Term	α	p = Number of Treatments								
		2	3	4	5	6	7	8	9	10
10	0.05	3.15	3.88	4.33	4.65	4.91	5.12	5.30	5.46	5.60
	0.01	4.48	5.27	5.77	6.14	6.43	6.67	6.87	7.05	7.21
11	0.05	3.11	3.82	4.26	4.57	4.82	5.03	5.20	5.35	5.49
	0.01	4.39	5.15	5.62	5.97	6.25	6.48	6.67	6.84	6.99
12	0.05	3.08	3.77	4.20	4.51	4.75	4.95	5.12	5.27	5.39
	0.01	4.32	5.05	5.50	5.84	6.10	6.32	6.51	6.67	6.81
13	0.05	3.06	3.73	4.15	4.45	4.69	4.88	5.05	5.19	5.32
	0.01	4.26	4.96	5.40	5.73	5.98	6.19	6.37	6.53	6.67
14	0.05	3.03	3.70	4.11	4.41	4.64	4.83	4.99	5.13	5.25
	0.01	4.21	4.89	5.32	5.63	5.88	6.08	6.26	6.41	6.54
15	0.05	3.01	3.67	4.08	4.37	4.59	4.78	4.94	5.08	5.20
	0.01	4.17	4.84	5.25	5.56	5.80	5.99	6.16	6.31	6.44
16	0.05	3.00	3.65	4.05	4.33	4.56	4.74	4.90	5.03	5.15
	0.01	4.13	4.79	5.19	5.49	5.72	5.92	6.08	6.22	6.35
17	0.05	2.98	3.63	4.02	4.30	4.52	4.70	4.86	4.99	5.11
	0.01	4.10	4.74	5.14	5.43	5.66	5.85	6.01	6.15	6.27
18	0.05	2.97	3.61	4.00	4.28	4.49	4.67	4.82	4.96	5.07
	0.01	4.07	4.70	5.09	5.38	5.60	5.79	5.94	6.08	6.20
19	0.05	2.96	3.59	3.98	4.25	4.47	4.65	4.79	4.92	5.04
	0.01	4.05	4.67	5.05	5.33	5.55	5.73	5.89	6.02	6.14
20	0.05	2.95	3.58	3.96	4.23	4.45	4.62	4.77	4.90	5.01
	0.01	4.02	4.64	5.02	5.29	5.51	5.69	5.84	5.97	6.09
24	0.05	2.92	3.53	3.90	4.17	4.37	4.54	4.68	4.81	4.92
	0.01	3.96	4.55	4.91	5.17	5.37	5.54	5.69	5.81	5.92
30	0.05	2.89	3.49	3.84	4.10	4.30	4.46	4.60	4.72	4.83
	0.01	3.89	4.45	4.80	5.05	5.24	5.40	5.54	5.65	5.76
40	0.05	2.86	3.44	3.79	4.04	4.23	4.39	4.52	4.63	4.74
	0.01	3.82	4.37	4.70	4.93	5.11	5.27	5.39	5.50	5.60

For our greenhouse example we get: $w = q_{.05(4,20)}\sqrt{(3.052/6)} = 3.96(0.7132) = 2.824$

Step 2: Rank the means, calculate differences

For the greenhouse example, we rank the means as:

29.20	28.6	25.87	21.00
-------	------	-------	-------

Start with the largest and second-largest means and calculate the difference, $29.20 - 28.60 = 0.60$, which is *less* than our w of 2.824, so we indicate there is no significant difference between these two means by placing the letter "a" under each:

29.20	28.6	25.87	21.00
a	a		

Then calculate the difference between the largest and third-largest means, $29.20 - 25.87 = 3.33$, which exceeds the critical w of 2.824, so we can label these with a "b" to show this difference is significant:

29.20	28.6	25.87	21.00
a	a	b	

Now we have to consider whether or not the second-largest and third-largest differ significantly. This is a step that sets up a back-and-forth process. Here $28.6 - 25.87 = 2.73$, less than the critical w of 2.824, so these two means do not differ significantly. We need to add a factor of "b" to show this:

29.20	28.6	25.87	21.00
a	ab	b	

Continuing down the line, we now calculate the next difference: $28.60 - 21.00 = 7.60$, exceeding the critical w , so we now add a "c":

29.20	28.6	25.87	21.00
a	ab	b	c

Again, we need to go back and check to see if the third-largest also differs from the smallest: $25.87 - 21.00 = 4.87$, which it does. So we are done.

These letters can be added to figures summarizing the results of the ANOVA.

The Tukey procedure explained above is valid only with equal sample sizes for each treatment level. In the presence of unequal sample sizes, more appropriate is the Tukey-Cramer Method, which calculates the standard deviation for each pairwise comparison separately. This method is available in SAS, R, and most other statistical softwares.

This page titled [2.3: Tukey Test for Pairwise Mean Comparisons](#) is shared under a [CC BY-NC 4.0](#) license and was authored, remixed, and/or curated by [Penn State's Department of Statistics](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.