

Unit 7: Topics in inference

1. ANOVA

Stat 140 - 02

Mount Holyoke College

Dr. Shan Shan

Slides posted at <http://sshanshans.github.io/stat140>

1. ANOVA: Analysis of Variance

Hypothesis:

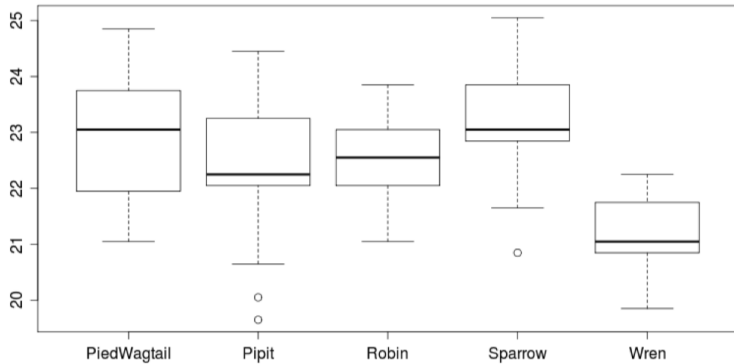
$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

Analysis of Variance (ANOVA) is a statistical method used to test differences between two or more means. It may seem odd that the technique is called “Analysis of Variance” rather than “Analysis of Means”. As you will see, the name is appropriate because inferences about means are made by analyzing variance.

- ▶ Cuckoo birds lay their eggs in the nests of other birds
- ▶ When the cuckoo baby hatches, it kicks out all the original eggs/babies
- ▶ If the cuckoo is lucky, the mother will raise the cuckoo as if it were her own
- ▶ Do cuckoo birds found in nests of different species differ in size?



Length of Cuckoo Eggs



- k = number of groups
- n_j = number of units in group j
- n = overall number of units
 $= n_1 + n_2 + \dots + n_k$

Bird	Sample Mean	Sample SD	Sample Size
Pied Wagtail	22.90	1.07	15
Pipit	22.50	0.97	60
Robin	22.58	0.68	16
Sparrow	23.12	1.07	14
Wren	21.13	0.74	15
Overall	22.46	1.07	120

- $k = 5$
- $n_1 = 15, n_2 = 60, n_3 = 16, n_4 = 14, n_5 = 15$
- $n = 120$

Null hypothesis:

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$$

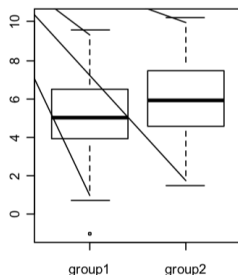
Poll question

Which of the following is a correct statement of the alternative hypothesis?

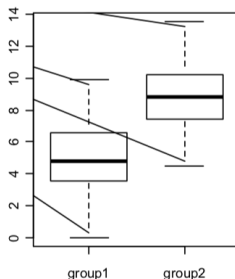
1. For any two groups, no two group means are the same.
2. There are at least two groups that have different group means from each other.

Whether or not two means are significantly different depends on

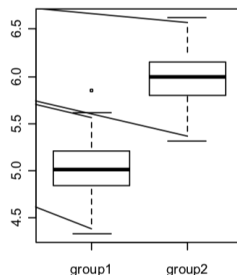
- ▶ How far apart the means are
- ▶ How much variability there is within each group



$$\begin{aligned}\bar{X}_1 &= 5 \\ \bar{X}_2 &= 6 \\ s_1 &= s_2 = 2\end{aligned}$$



$$\begin{aligned}\bar{X}_1 &= 5 \\ \bar{X}_2 &= 9 \\ s_1 &= s_2 = 2\end{aligned}$$



$$\begin{aligned}\bar{X}_1 &= 5 \\ \bar{X}_2 &= 6 \\ s_1 &= s_2 = 0.2\end{aligned}$$

Poll question

If the groups are actually different, then

- a the variability between groups should be higher than the variability within groups
- b the variability within groups should be higher than the variability between groups

Analysis of Variance(ANOVA) compares the variability between groups to the variability within groups

$$\text{Total Variability} = \text{Variability Between Groups} + \text{Variability Within Groups}$$

The F-statistic is a ratio of the average variability between groups to the average variability within groups

$$\text{F-statistic} = \frac{\text{average variability between groups}}{\text{average variability within groups}}$$

Poll question

If there really is a difference between the groups, we would expect the F-statistic to be

- a Higher than we would observe by random chance
- b Lower than we would observe by random chance

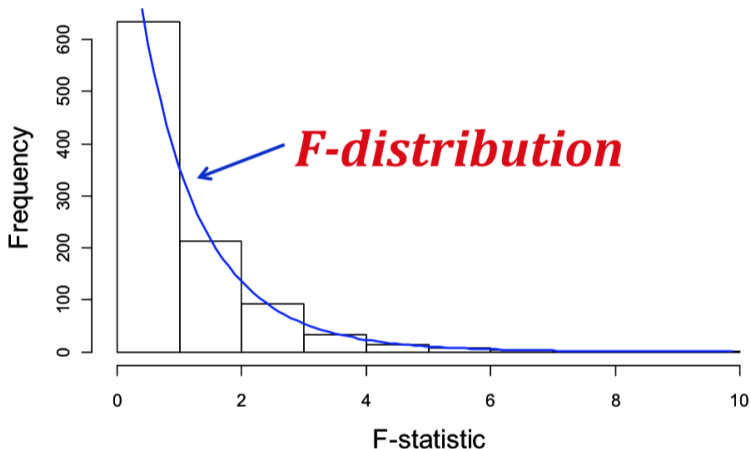
Poll question

If there really is a difference between the groups, we would expect the F-statistic to be

- a Higher than we would observe by random chance
- b Lower than we would observe by random chance

We now have a test statistic. What else do we need to perform the hypothesis test?

A distribution of the test statistic assuming H_0 is true



The F-distribution has two degrees of freedom, one for the numerator of the ratio $(k-1)$ and one for the denominator $(n-k)$

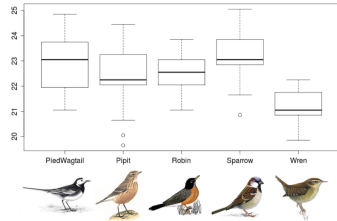
If the following conditions hold,

1. Sample sizes in each group are large (each $n_j \geq 30$) OR the data are relatively normally distributed
2. Variability is similar in all groups
 - As a rough rule of thumb, this assumption is violated if the standard deviation of one group is more than double the standard deviation of another group

Under the null hypothesis, the F-statistic follows an F-distribution

Bird	Sample Mean	Sample SD	Sample Size
Pied Wagtail	22.90	1.07	15
Pipit	22.50	0.97	60
Robin	22.58	0.68	16
Sparrow	23.12	1.07	14
Wren	21.13	0.74	15
Overall	22.46	1.07	120

- $k = 5$
- $n_1 = 15, n_2 = 60, n_3 = 16, n_4 = 14, n_5 = 15$
- $n = 120$



Poll question

Can we use the F-distribution to calculate the p-value for the Cuckoo bird eggs?

- a Yes
- b No
- c Need more information

How likely it is to observe a test statistic as extreme (or more extreme) under the F -distribution?

Use the 'aov()' function

```
aov.out = aov(group ~ lengths, data=cuckoo)
summary(aov.out)
```

If $p\text{-value} < \text{significance level}$, reject H_0 ;
Otherwise, do not reject H_0 .

Use t-tests and the Bonferroni correction

- ▶ If the ANOVA yields a significant results, next natural question is: “Which means are different?”
- ▶ Use t-tests comparing each pair of means to each other
- ▶ Compare resulting p-values to a modified significance level

$$\alpha^* = \frac{\alpha}{K}$$

where $K = k(k - 1)/2$ is the total number of pairwise tests