

Introduction

In order to avoid further revisiting of any second exam traumas, we'll abandon the espresso dataset as a motivating example.

Instead, we'll use data collected from the National Advanced Driving Simulator (NADS) for a study investigating the link between drug use and risky driving behavior.

Figure 1: An Example of Risky Driving Behavior



Solution

We first obtain a 95% confidence interval for the difference in mean $\log(Distance)$ between the two groups:

$$\overline{\log(D_{ND})} - \overline{\log(D_{THC})} \pm t_{crit} \sqrt{\frac{S_{ND}^2}{n_{ND}} + \frac{S_{THC}^2}{n_{THC}}} = (-0.151, 0.318)$$

Exponentiating the endpoints of this interval, we obtain the 95% CI for the mean relative increase in following distance of No Drug and THC users:

$$(\exp(-0.151), \exp(0.318)) = (0.86, 1.37)$$

The above result indicates that the No Drug following distance plausibly ranges from 14% shorter to 37% longer than the THC following distance.

The test statistic on the log scale is 0.71 with a p-value of 0.478, and, on the original scale, the test statistic is 0.39 with a p-value of 0.70.

While both would lead us to fail to reject the null hypothesis, the test on the log-scale is more powerful since it better meets the t-test's normality assumptions.



Weren't We Supposed to Talk About ANOVA?

We began this lecture by teasing an analytic approach called ANOVA, and subsequently discussed the tangential topics of outliers and log-transformations.

With these important concepts behind us, we'll next speak broadly about making comparisons across multiple groups, and then lead into a discussion of statistical modeling and ANOVA.



Solution

- 1) ALC vs NODRUG, p-value = 0.5102
- 2) ALC vs MDMA, p-value = 0.00417
- 3) ALC vs THC, p-value = 0.8959
- 4) THC vs NODRUG, p-value = 0.4782
- 5) THC vs MDMA, p-value = 0.01383
- 6) MDMA vs NODRUG, p-value = 0.00216

Using the unadjusted threshold, there are three significant results: 2), 5), and 6). However, we expect our overall type I error rate to be higher than 0.05.

Using the Bonferroni-adjusted threshold ($0.05/6 = 0.0083$), we find only two significant results: 2) and 6).

We would then conclude that the following distances are different between the ALC and MDMA groups, as well as between the MDMA and NODRUG groups.



What About Power?

While using a Bonferroni adjustment effectively controls our overall type I error rate, it comes at the cost of statistical power.

Applying a Bonferroni adjustment decreases the significance level, α , of each individual test. As a result, the amount of evidence necessary to reject each null hypothesis increases, making the number of rejections fewer than if using an unadjusted α .

With fewer rejections, we can then expect a decrease in power (and increase in type II error rate).

If we wanted to avoid this tradeoff and still control our type I error rate, we could turn to option (2) and use a single, joint test (i.e. ANOVA) of the hypothesis:

$$H_0 : \mu_{ND} = \mu_{THC} = \mu_{ALC} = \mu_{MDMA}$$



ANOVA

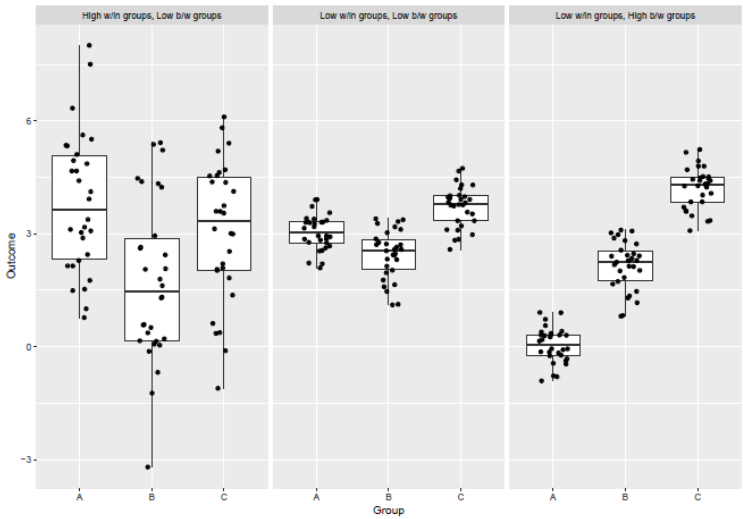
As the name implies (i.e. Analysis of Variance), in ANOVA we analyze the total observed variability in our outcome data in order to appropriately partition it between two sources:

- 1) Variability arising due to differences between groups (i.e. between-group variability)
- 2) Variability arising due to differences within groups (i.e. within-group variability)

Intuitively, if there was more between-group variability than within-group variability, it would be reasonable to conclude that there are significant group-level differences in our outcome.



Partitioning Variability



Statistical Modeling

To better frame the idea of partitioning variability, it is important that we first discuss **statistical modeling**.

A model is a simplified characterization of a certain process or relationship. Good models are those which, while maintaining some degree of simplicity, accurately describe or explain the phenomenon of interest.

Statistical models are models in every sense of the word, but their "quality" is measured by their ability to "explain" the variability in a certain outcome variable.

It is generally impossible for models to explain all of the variability in an outcome, but some models are better than others in this endeavor.

As an example, consider the following two models:

- 1) Using the height of a child's parents to predict the child's adult height
- 2) Using the child's weight at birth to predict the child's adult height.

Clearly, we would expect the first model to do a better job in predicting, or explaining the variability in, the child's adult height.



Statistical Models

The simplest statistical model, often referred to as the **null model**, is one which posits that all outcome variability is "unexplainable" and should therefore be modeled using a single mean.

In the case of the NADS study, the null model would correspond to using the mean following distance (across all groups) as the prediction for everyone in the study.

More complex statistical models are those which use one or more *explanatory* variables to explain the variability in an outcome of interest.

For the NADS study, an example of a more complex statistical model would involve using the mean following distance for a study group (NODRUG, ALC, THC, or MDMA) as the prediction for individuals in that group.



Total Variability

The quality of a model is measured by its ability to "explain" the variability in a certain outcome variable.

We often summarize the total variability in an outcome by the **Total Sum of Squares (SST)**:

$$SST = \sum_i (y_i - \bar{y})^2$$

You may recognize that this is simply the sum of the **residuals**, $r_i = y_i - \bar{y}$.

We can then assess the quality (or fit) of a model by the share of SST it explains. The greater share, or proportion, of SST that is accounted for by a model, the better that that model is.



Sum of Squares Error

Since models generally are not perfect, it is often of interest to quantify the amount by which their predictions fall short.

We often use the **Sum of Squares Error**, or **SSE**, as a measure of model error. The larger the SSE, the worse a model is.

The way in which the SSE is computed is dependent on the model that is fit.

As an example, for the NADS study, the SSE for the null model is computed:

$$\sum_i (y_i - \bar{y})^2$$

In contrast, if we wanted to use the study groups to model our outcome, the SSE would be computed:

$$\sum_i (y_i - \bar{y}_i)^2$$

where \bar{y}_i is the appropriate group-specific mean.



Model Fit

Using both the SSE and SST, we can quantify the fit of a model (i.e. proportion of explained variability) through the **coefficient of determination**, R^2 :

$$R^2 = \frac{SST - SSE}{SST}$$

Under the null model, $SST = SSE$, and so R^2 is 0.

For the tailgating data, under the model where each group gets its own mean, $R^2 = 0.055$ indicating that the model explains 5.5% of the total variability.

Increasing the complexity of a model will always lower the SSE and increase R^2 . Despite this, care must be taken not to **overfit** your model by adding unnecessary complexity.

Rather, we should test whether introducing complexity lowers the SSE by more than what we'd expect to see by random chance.



ANOVA

ANOVA is simply a special type of statistical model in which a single categorical variable is used to predict a quantitative outcome.

In order to test whether the drop in SSE resulting from including a single categorical variable is greater than what we'd expect to see by random chance, we use the test statistic:

$$F = \frac{(SST - SSE)/(d_1 - d_0)}{\text{Std. Error}}$$

where d_1 and d_0 refer to the number of parameters in the model being considered (i.e. the model using a single categorical variable as a predictor) and the null model.

For the NADS study, $d_0 = 1$ (the single overall mean) and $d_1 = 4$ (each group's mean).



ANOVA

For tests we've learned in the past that involve quantitative variables, we've seen that the standard errors involve some measure of variability divided by the sample size (e.g. $\frac{s}{\sqrt{n}}$).

In the ANOVA setting:

$$\text{Std. Error} = \frac{SSE}{n - d_1}$$

Using this standard error, the F statistic can be expressed:

$$F = \frac{(SST - SSE)/(d_1 - d_0)}{SSE/(n - d_1)}$$



F-Test and Variability

SST is the total sum of squares which quantifies the total variability in the outcome, y . $SST = \sum_i (y_i - \bar{y})^2$.

SSE is the outcome variability that remains unexplained after implementing your model. Under the null model, $SSE = SST$. Under the ANOVA model, $SSE = \sum_i (y_i - \bar{y}_i)^2$.

By subtraction, we can determine how much variability is being explained by the inclusion of our categorical variable:

$$SSG = SST - SSE$$

where SSG , the **Sum of Squares Groups**, quantifies the amount of variability explained using the categorical variable groups in our model.



F-Test and Variability

Using SSG, we can express the F-statistic as:

$$F = \frac{SSG/(d_1 - d_0)}{SSE/(n - d_1)}$$

Sums of squares divided by their degrees of freedom are often called **mean squares**, and allows us to express the F-statistic as:

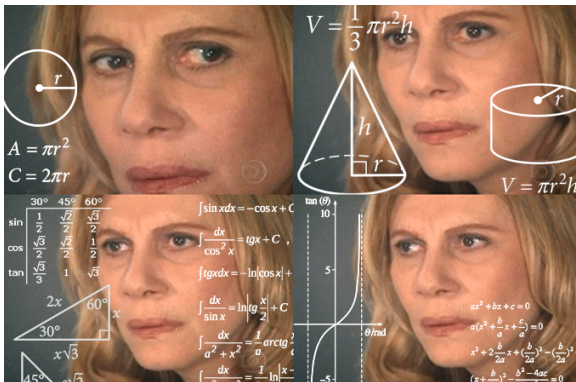
$$F = \frac{MSG}{MSE}$$

where MSG is the mean square of groups and MSE is the mean square of errors.



Math Overload

At this point in the lecture, and with all the math thrown at you all, I'm seeing a lot of this:



Math Overload

This is completely understandable, and the good news is that I will not expect you to memorize these formulas or compute any of these things by hand.

However you will be expected to understand and interpret an **ANOVA table**, which is a common piece of software output for ANOVA.

Then general form of these tables is shown below:

| Source | DF | SS | MS | F-Value | P-Value |
|--------|-------------|-------|-------|-----------|--------------------------|
| Group | $d_1 - d_0$ | SSG | MSG | MSG/MSE | Use $F_{d_1-d_0, n-d_1}$ |
| Error | $n - d_1$ | SSE | MSE | | |
| Total | $n - d_0$ | SST | | | |

$d_0 = 1$, the null model has one parameter, a single overall mean

$d_1 = k$, the alternative model has k parameters, a different mean for each group.



Practice

With your groups, complete the following ANOVA table (assuming this is a typical ANOVA test with $d_0 = 1$):

| Source | DF | SS | MS | F-Value | P-Value |
|--------|----|-----|----|---------|---------|
| Group | 4 | 200 | ? | ? | ? |
| Error | ? | 440 | ? | | |
| Total | 59 | ? | | | |



Solution

In this example, $d_1 = k = 5$ and $n = 60$ so:

| Source | DF | SS | MS | F-Value | P-Value |
|--------|----|-----|----|---------|---------|
| Group | 4 | 200 | 50 | 6.25 | 0.0003 |
| Error | 55 | 440 | 8 | | |
| Total | 59 | 640 | | | |

The p-value is found using the right-tail area beyond 6.25 of an F distribution with (4,55) degrees of freedom.



Minitab Practice

With your groups, use Minitab to analyze the Tailgating data with ANOVA (Stat – > ANOVA – > One-way). Use $\log(\textit{Distance})$ as your outcome variable. Be sure to report:

- 1) Your null and alternative hypotheses
- 2) Your test-statistic
- 3) Your p-value and a one sentence conclusion



Solution

- 1) $H_0 : \mu_{ND} = \mu_{THC} = \mu_{ALC} = \mu_{MDMA}$; $H_A : \mu_i \neq \mu_j$ for at least one pair.
- 2) $F = 2.23$
- 3) Using a $F_{3,115}$ distribution, we obtain a p-value of 0.088. There is borderline evidence that drug use is predictive of following distance. It appears that the MDMA group is most different in that the group generally has shorter following distances than the rest.



Inference After ANOVA

The results of an ANOVA test only tells us whether a difference across groups exists, and not which specific groups are different.

Because of this, ANOVA is often followed by a few pairwise comparisons in order to investigate which groups differ.

In Minitab we can do this using **Tukey's honest significant difference (HSD) test** (sometimes called Tukey's range test).

Similar to the Bonferroni adjustment, Tukey's HSD controls the family type I error rate for all possible pairwise comparisons (so we don't need to worry about an inflated overall type I error rate).



More on ANOVA and Modeling

In ANOVA, we use a single categorical variable to predict a quantitative outcome variable.

The ANOVA test will be statistically significant only if the categorical variable improves prediction beyond what could be attributed to random chance.

The ANOVA model is just one type of model in the vast array of statistical models. If we were to cover all statistical models, you'd be halfway towards getting a PhD in statistics!

We'll instead restrict the last of our lectures to regression models, of which ANOVA is a special case.

Aside from ANOVA, regression models include the simple regression models seen earlier in the course (i.e. single predictor) and multiple regression models with several predictors.



Wrap-Up

Right now, you should...

- Understand how ANOVA testing is related to statistical modeling
- Understand the partitioning of variability in ANOVA
- Construct and use an ANOVA table to draw conclusions
- Conduct appropriate follow-up analyses after ANOVA

These notes cover sections 8.1-8.2 of the textbook. Please read through these sections and their examples along with any links provided in this lecture.

