

MET CS 555 - Data Analysis and Visualization

Module-5: Anova and Regression, ANCOVA

Lecture - 10

Kia Teymourian

Boston University

Slides last compiled: 02/22/2019, Time: 01:28:30

Table of contents

1. Anova and Regression
2. One-Way Analysis of Covariance (ANCOVA)
3. Two-Way Analysis of Variance

Anova and Regression

One-Way Analysis of Variance and Regression

We can use the linear regression to conduct the same tests in the one-way ANOVA setting

The one-way ANOVA is the **same as performing a regression** where the explanatory variable in the model is a variable or variables that indicate group membership.

In order to represent a one-way ANOVA model in the regression framework, construction of dummy variables is required.

For a particular categorical variable with **k categories, k-1 dummy variables are needed.**

$$\begin{aligned}\text{group}_2 &= \begin{cases} 1, & \text{if observation is in group 2} \\ 0, & \text{otherwise} \end{cases} \\ \text{group}_3 &= \begin{cases} 1, & \text{if observation is in group 3} \\ 0, & \text{otherwise} \end{cases} \\ &\vdots \\ \text{group}_k &= \begin{cases} 1, & \text{if observation is in group } k \\ 0, & \text{otherwise} \end{cases}\end{aligned}$$

One-Way Analysis of Variance and Regression

For a particular categorical variable with **k categories**, **k-1 dummy variables are needed**.

Dummy variables are binary variables of the form:

The category omitted (group 1 in this case) is referred to as the reference group.

A dummy variable for the reference group is not needed as those in this group can be identified as having values of 0 for all of the defined dummy variables.

$$\begin{aligned}\text{group}_2 &= \begin{cases} 1, & \text{if observation is in group 2} \\ 0, & \text{otherwise} \end{cases} \\ \text{group}_3 &= \begin{cases} 1, & \text{if observation is in group 3} \\ 0, & \text{otherwise} \end{cases} \\ &\vdots \\ \text{group}_k &= \begin{cases} 1, & \text{if observation is in group } k \\ 0, & \text{otherwise} \end{cases}\end{aligned}$$

One-Way Analysis of Variance and Regression

The following table summarizes how dummy variables would be defined for a grouping variable with k categories.

If you knew the values of each of the dummy variables for a particular observation, you would then know what group the observation was a member of.

Group	Dummy Variables			
	group ₂	group ₃	...	group _k
1	0	0		0
2	1	0		0
3	0	1		0
...				
k	0	0		1

One-Way Analysis of Variance and Regression

The following model can be used to compare means between groups in the regression framework:

$$y = \beta_0 + \sum_{i=2}^k \beta_{i-1} \text{group}_i + e$$

where

- ▷ y is the response or dependent variable. $group_2, group_3, \dots, group_k$ are the dummy variables that represent membership in groups 2, 3, ... and k , respectively.
- ▷ β_0 is the intercept (the sample mean in the reference group [group 1 in this case]).
- ▷ β_1 is the mean difference between group 2 and the reference group (group 1 in this case).
- ▷ β_{k-1} is the mean difference between group k and the reference group.
- ▷ e is the random error which we assume is normally distributed with a mean of 0 and a variance of σ^2 .

One-Way Analysis of Variance and Regression

We can use this model to test the hypotheses of interest in the one-way ANOVA setting.

We can use this model to test the typical ANOVA null hypothesis

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

(All underlying population means are equal) against the alternative hypothesis

$$H_1 : \mu_i \neq \mu_j \text{ for some } i \text{ and } j$$

(At least two of the k underlying population means are different or not all of the underlying population means are the same/equal).

One-Way Analysis of Variance and Regression

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

(All underlying population means are equal) against the alternative hypothesis

$$H_1 : \mu_i \neq \mu_j \text{ for some } i \text{ and } j$$

(At least two of the k underlying population means are different or not all of the underlying population means are the same/equal).

Mathematically, these hypotheses are equivalent to the global F-test hypotheses in regression where we test the null hypothesis that all slope coefficients are equal to 0 .

$H_0 : \beta_1 = \beta_2 = \cdots = \beta_k = 0$ versus the alternative that at least one of the slope coefficients is different from zero

$$H_1 : \beta_i \neq 0 \text{ for at least one } i$$

Given that the **overall F-test for one-way ANOVA is equivalent to the global F-test for multiple linear regression (when dummy variables are used in the regression to represent group membership)**, the ANOVA table in each case is the same.

One-Way Analysis of Variance and Regression

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

$$H_1 : \mu_i \neq \mu_j \text{ for some } i \text{ and } j$$

Mathematically equivalent to

$H_0 : \beta_1 = \beta_2 = \cdots = \beta_k = 0$ versus the alternative that at least one of the slope coefficients is different from zero

$$H_1 : \beta_i \neq 0 \text{ for at least one } i$$

This also means that the **sum of squares between** and the **sum of squares within** from the **one-way ANOVA model** are equivalent to the regression **sum of squares** and the **residual sum of squares**, respectively.

R functions for One-way ANOVA using lm()

Create dummy variables and One-way ANOVA using lm() function

```
# Create dummy variables
> data$g0 <- ifelse(data$group=='Current heavy smoker', 1, 0)
> data$g1 <- ifelse(data$group=='Current light smoker', 1, 0)
> data$g2 <- ifelse(data$group=='Former smoker', 1, 0)
> data$g3 <- ifelse(data$group=='Never smoker', 1, 0)

# One-way ANOVA using lm() function
> m2 <- lm(data$SBP~data$g0+data$g1+data$g2, data=data)
> summary(m2)

> m3 <- lm(data$SBP~data$g1+data$g2+data$g3, data=data)
> summary(m3)

> m4 <- lm(data$SBP~data$g0+data$g2+data$g3, data=data)
> summary(m4)
```

The Golf Ball Example

In the golf ball example, the one-way ANOVA table was as follows:

	SS (Sum of Squares)	df (degrees of freedom)	MS (Mean Square)	F
Between	2583.3335	2	1291.6668	20.67
Within	750	12	62.5	
Total	3333.3335			

Let us we create dummy variables for 2 of the 3 brands, for example, for the Nike and Callaway brands, and ran a regression predicting distance from these two dummy variables.

The regression model in this case that is equivalent to the above one-way ANOVA model is given by

$$y = \beta_0 + \beta_{\text{Nike}} \text{group}_{\text{Nike}} + \beta_{\text{Callaway}} \text{group}_{\text{Callaway}} + e$$

where $group_i$ is a dummy variable indicating whether or not the observation is of brand i .

The golf example: R commands

```
# creating dummy variables
> golf$g0 <- ifelse(brand=='Callaway', 1, 0)
> golf$g1 <- ifelse(brand=='Nike', 1, 0)
> golf$g2 <- ifelse(brand=='Titleist', 1, 0)

# Create the model
> m <- lm(dist ~ g1 + g2, data=golf)
> summary(m)

# pass the model to the anova function
> anova(m)
```

```
Call:
lm(formula = dist ~ g1 + g2, data = golf)

Residuals:
    Min       1Q   Median       3Q      Max
   -10      -5         0         5      10

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   285.000     3.536   80.61 < 2e-16 ***
g1             -25.000     5.000   -5.00 0.000309 ***
g2              5.000     5.000    1.00 0.337049
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.906 on 12 degrees of freedom
Multiple R-squared:  0.775,    Adjusted R-squared:  0.7375
F-statistic: 20.67 on 2 and 12 DF,  p-value: 0.0001297
```

The golf example: R commands

```
# creating dummy variables
> golf$g0 <- ifelse(brand=='Callaway', 1, 0)
> golf$g1 <- ifelse(brand=='Nike', 1, 0)
> golf$g2 <- ifelse(brand=='Titleist', 1, 0)
# Create the model
> m <- lm(dist~g1+g2, data=golf)
> summary(m)
# pass the model to the anova function
> anova(m)
```

Analysis of Variance Table

Response: dist

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
g1	1	2520.8	2520.8	40.333	3.66e-05 ***
g2	1	62.5	62.5	1.000	0.337
Residuals	12	750.0	62.5		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
, |

Pairwise comparisons can also be performed in the regression

framework. The most simple is the test of the underlying mean from group i to the reference group (group 1 as outlined above).

In the one-way ANOVA framework, the test of the null hypothesis $\mu_i = \mu_1$ would be accomplished through the use of the t statistic:

$$t = \frac{\bar{x}_i - \bar{x}_1}{\sqrt{s_p^2 \left(\frac{1}{n_i} + \frac{1}{n_1} \right)}}$$

which follows a t-distribution with $n-k$ degrees of freedom under H_0 (k here is the number of groups).

This test is the same as the t-test in the regression setting which tests the null hypothesis $\beta_i = 0$ after controlling for the other independent variables in the model.

Pairwise comparisons - t-test

$$t = \frac{\bar{x}_i - \bar{x}_1}{\sqrt{s_p^2 \left(\frac{1}{n_i} + \frac{1}{n_1} \right)}}$$

The t-statistic here is the same as the t-statistic for multiple regression:

$$t = \frac{\hat{\beta}_i}{\text{SE}_{\hat{\beta}_i}}$$

which follows a t-distribution with $n - k$ degrees of freedom under H_0

(k here represents the number of groups and not the number of variables in the model as we had defined it in the regression context).

Pairwise comparisons - t-test

It is also possible to perform other **pairwise comparisons** in the regression setting (for example, to test if there is a difference between groups i and j where neither i nor j represent the reference group), it **requires testing of null hypotheses of the form $\beta_i = \beta_j$ which are not standard** output from the regression model in most statistical packages and thus requires additional steps in order to calculate the test statistic for such comparisons.

- ▷ The same comparisons are relatively easy to make from the one-way ANOVA setting.
- ▷ If one were interested in using regression to make these comparisons, the **easiest way to do it is to run more than one regression analysis and change the reference group each time to a different group.**
- ▷ Obviously, conducting one-way ANOVAs gives this output standardly as such is often easier to implement in statistical software packages.

The golf example

For the golf ball distance by brand example, the below table summarizes the results of the **pairwise comparisons**. The group means were 285, 260 and 290 for Callaway, Nike, and Titleist, respectively.

Comparison	Mean Difference	<i>t</i>	df	Unadjusted p-value
Titleist versus Callaway	5	1.0	12	0.34
Titleist versus Nike	30	6.0	12	0.000062
Callaway versus Nike	25	5.0	12	0.00031

The summary of the output of the **regression model** below are show in the table:

$$\hat{y} = \beta_0 + \beta_{\text{Nike}} \text{group}_{\text{Nike}} + \beta_{\text{Callaway}} \text{group}_{\text{Callaway}}$$

	Estimate	SE	<i>t</i>	$Pr(> t)$
Intercept	290	3.536	82.02	< 0.0001
Nike versus Titleist	-30	5	-6	0.000062
Callaway versus Titleist	-5	5	-1	0.34

One-Way Analysis of Covariance (ANCOVA)

One-Way Analysis of Covariance (ANCOVA)

Sometimes we have **other variables or factors that have an effect on the dependent variables** that we are analyzing. But these variables are not actually independent variables.

You want to control the effect of other variables in your analysis.

For example: When you study the effect of different diet types on weight loss other factors like age or level of training can have an effect. You do not want to consider them as a separate independent variables, you just want to check if these variables have an effect on your analysis and control it.

We have multiple variables that have an effect on a single numeric outcome.

This is where we use **analysis of covariance**.

One-Way Analysis of Covariance (ANCOVA)

The other variables are variables that are called **“covariates”**.

These variables provide additional variations in the subject response value. We want to control these variables.

Covariates are also named “Intervening” or “Confounding Variables”.

Ancova is an extension of Anova and can be seen as a combination of Anova and Regression.

What ANCOVA does it simply takes out the effects of covariates out of Anova analysis.

One-Way Analysis of Covariance (ANCOVA)

- ▶ ANCOVA is a general linear model that is a **combination of ANOVA and regression** when you have some **categorical factors** and some **quantitative continuous variables**.
- ▶ The continuous variables (on which to perform regression) are called **“covariates”**
- ▶ Often these covariates are not necessarily of primary interest, but still their inclusion in the model will help explain more of the response, and hence reduce the error variance.
- ▶ ANCOVA evaluates whether population means of a response variable are equal across levels of a categorical explanatory variable, while statistically controlling for the effects of other continuous variables that are not of primary interest, known as covariates.
- ▶ Mathematically, ANCOVA decomposes the variance in the response variable into variance explained by the covariates, variance explained by the categorical explanatory variables, and residual variance.

One-Way Analysis of Covariance (ANCOVA)

The one-way ANCOVA model can be expressed as:

$$y = \beta_0 + \sum_{i=2}^k \beta_{i-1} \text{group}_i + \sum_{i=1}^j \beta_{k+i-1} x_i + e$$

where

- ▷ $\text{group}_2, \text{group}_3, \dots, \text{group}_k$ are the dummy variables that represent membership in groups 2, 3, ... and k , respectively.
- ▷ β_0 is the intercept (the sample mean in the reference group [group 1 in this case]).
- ▷ β_1 is the mean difference between group 2 and the reference group (group 1 in this case).
- ▷ β_{k-1} is the mean difference between group k and the reference group.
- ▷ β_k is the expected change in y for each one unit change in x_1 , across all groups and after adjusting for the other covariates x_2, x_3, \dots, x_j .
- ▷ β_{k+j-1} is the expected change in y for each one unit change in x_j , across all groups and after adjusting for the other covariates x_1, x_2, \dots, x_{j-1} .

One-Way ANCOVA - Inferencing

In one-way ANCOVA framework, we are mainly interested in the overall model results as well as the results for the grouping factor (after adjusting for other variables in the model).

We can test the overall null hypothesis

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_{k+j} = 0$$

(none of the variables in the model is predictive of the dependent variable y)

versus the alternative hypothesis

$$H_1 : \text{there is at least one } \beta_i \neq 0$$

(at least one of the variables is predictive of the dependent variable y)

One-Way ANCOVA - Inferencing

If the global null hypothesis is rejected, then we can test if the underlying population means are different across groups after controlling for the other variables in the model.

$H_0 : \beta_1 = \beta_2 = \dots = \beta_{k-1} = 0, \beta_k \neq 0, \dots, \beta_{k+j} \neq 0$ (all underlying population means are equal after controlling for x_1, \dots, x_j)

versus the alternative hypothesis

$H_1 : \text{At least two of the } k \text{ underlying population means are different}$
or not all of the underlying population means are the same/equal after controlling for x_1, \dots, x_j

One-Way ANCOVA - least-squares means or LS means

In the ANCOVA setting, if there is a difference between groups after adjusting for the other covariates in the model, we follow up with a comparison for the adjusted group means (instead of a comparison of the sample means).

The adjusted means are often called least-squares means or LS means.

- ▷ The **LS means** represent the **mean value for each group** that is adjusted for the other **covariates included in the model**.
- ▷ The **LS mean** in each group is calculated using the **least-squares regression equation** using the mean values of the covariates in the model.
- ▷ The general procedure and interpretation of the **one-way ANCOVA model** is similar to the procedure for the **one-way ANOVA**.
- ▷ **The main difference is that in the ANCOVA setting, the inferences made are based on comparisons made after adjusting for or controlling for the covariates included in the model.**

One-Way ANCOVA

```
# Don't use aov() function as it will not produce expected output
# Need to use Anova() function from package car to get Type
    III sums of
square

> install.packages("car")
> library(car)
> Anova(lm(data$response ~ data$group + data$covariate), type=3)
```

car: Companion to Applied Regression

Functions and Datasets to Accompany J. Fox and S. Weisberg,
An R Companion to Applied Regression, Second Edition, Sage, 2011.
<https://cran.r-project.org/web/packages/car/index.html>

An example: One-Way ANCOVA

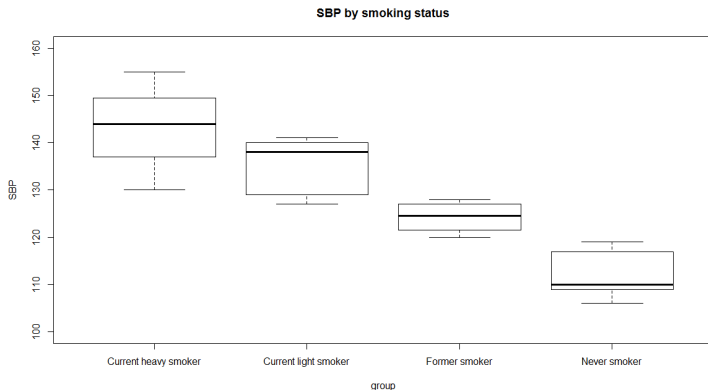
A random sample of current light smokers, current heavy smokers, former smokers, and those who have never smoked was taken to determine if mean systolic blood pressure (SBP) differs across smoking status categories.

Given that it is known that SBP increases with advancing age and given the fact that smoking preference are known to differ by age, we'd want to conduct an ANCOVA so that we can be sure that differences we see between smoking categories aren't due to differences in age.

It is preferable in this case to **use the ANCOVA methodology so that we can adjust for age.**

An example: One-Way ANCOVA

```
> data <- read.csv("smoking_SBP.csv")
> data
> aggregate(data$SBP, by=list(data$group), summary)
> boxplot(data$SBP~data$group, data=data, main="SBP by smoking status",
          xlab="group", ylab="SBP", ylim=c(100, 160))
```



An example: One-Way ANCOVA

First, ran a one-way ANOVA (without adjustment for age).

The global F-test showed that mean SBP differed by smoking category ($F=21.49$ on 3 and 15 degrees of freedom, $p < 0.001$).

```
> m<- aov(data$SBP~data$group, data=data)
> summary(m)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
data\$group	3	2786.2	928.7	21.49	1.1e-05 ***
Residuals	15	648.3	43.2		

```
---
```

Signif. codes:	0	***	0.001	**	0.01	*	0.05	.	0.1
----------------	---	-----	-------	----	------	---	------	---	-----

```
> qf(.95, df1=3, df2=15)
[1] 3.287382
```

An example: One-Way ANCOVA

After adjusting for multiple comparisons using Tukeys methodology, all pairwise comparisons except for (a) between light and heavy current smokers and (b) between light and former smokers were significant at the $\alpha=0.05$ level.

```
> is.factor(data$grpnum)
> fnum = factor(data$grpnum)
> fnum
> m1<- aov(data$SBP~fnum, data=data)
> TukeyHSD(m1)
```

```
Tukey multiple comparisons of means
 95% family-wise confidence level
```

```
Fit: aov(formula = data$SBP ~ fnum, data = data)
```

\$fnum		diff	lwr	upr	p adj
1-0	-8.25000	-20.96090	4.4608964	0.2811214	
2-0	-19.00000	-32.39846	-5.6015387	0.0047930	
3-0	-31.41667	-43.64773	-19.1856009	0.0000119	
2-1	-10.75000	-23.46090	1.9608964	0.1122314	
3-1	-23.16667	-34.64042	-11.6929099	0.0001780	
3-2	-12.41667	-24.64773	-0.1856009	0.0460500	

An example: One-Way ANCOVA

However, after adjusting for age using an ANCOVA model, the differences seen in the one-way ANOVA setting were attenuated and the F-test for the **effect of smoking status was no longer significant** ($F=1.774$ on 3 and 14 degrees of freedom, $p=0.1982$).

The least square means (adjusted for age) were 129.87, 144.01, 144.01, and 155.11 for the heavy, light, former and never smokers, respectively.

As such, the differences that we saw in the one-way ANOVA model were due to age differences across the smoking groups as opposed to true differences in SBP attributable only to smoking status.

```
> is.factor(data$grpnum)
> fnum = factor(data$grpnum)
> fnum
> m1<- aov(data$SBP~fnum, data=data)
> TukeyHSD(m1)
```


One-Way ANCOVA

```
#Re-run ANOVA adjusting for Age
> library(car)
> Anova(lm(data$SBP~data$group+data$age), type=3)

# Least square means
# install.packages("emmeans")
library(emmeans)
my.model<-lm(SBP~group+age, data = data)
emm_options(contrasts=c("contr.treatment", "contr.poly"))
emmeans(my.model, specs = "group")

# no p value adjustment
emmeans(my.model, specs = "group" , contr = "pairwise", adjust="none")

# P value adjustment: tukey method
emmeans(my.model, specs = "group" , contr = "pairwise", adjust="tukey")

# P value adjustment: bonferroni method for 6 tests
emmeans(my.model, specs = "group" , contr = "pairwise", adjust="
  bonferroni")
```

One-Way ANCOVA

```
# Generate least square means (covariate adjusted means) and comparisons

> install.packages('lsmeans')
> library(lsmeans)
> options(contrasts=c("contr.treatment", "contr.poly"))
> lsmeans(lm(data$response~data$group+data$covariate), pairwise~data$
  group, adjust=[method])

# Note: method =   tukey   ,   scheffe   ,   sidak", "bonferroni",
  dunnett",   mvt",
"none")
```

lsmeans: Least-Squares Means

Obtain least-squares means for many linear, generalized linear, and mixed models. Compute contrasts or linear functions of least-squares means, and comparisons of slopes. Plots and compact letter displays.

<https://cran.r-project.org/web/packages/lsmeans/vignettes/using-lsmeans.pdf>

Two-Way Analysis of Variance

Two-Way Analysis of Variance

Sometimes we are interested in comparing the **mean response across groups** when there is **more than one factor** to be considered.

When there are **two factors**(each with two or more levels) that we are interested in, we use the **two-way ANOVA methodology** instead of the one-way ANOVA methodology.

The **goal of a two-way ANOVA is to look at the effects of each factor after controlling for the effects of the other factor.**

Two-Way Analysis of Variance

Within the two-way ANOVA framework, we seek to test the following:

- ▷ **Whether or not the first factor impacts the mean outcome** after controlling for the **second factor**.

Here, we test H_0 : All underlying population means are equal across levels of the first factor, after controlling for the second factor.

- ▷ **Whether or not the second factor impacts the mean outcome** after controlling for the first factor.

Here, we test H_0 : All underlying population means are equal across levels of the second factor, after controlling for the **first factor**.

Two-Way Analysis of Variance

The alternative hypothesis in each case is that the **underlying populations means are not equal across levels of the factor tested after controlling for the other**. At least two of the underlying group means are different across the factor tested, after controlling for the other factor.

Consider the following notation:

r = No. of levels of factor A.

c = No. of levels of factor B.

μ_{ij} = population mean for those in the *ith* level of factor **A** and the *jth* level of factor **B**.

$\mu_{i.}$ = population mean for those in the *ith* level of factor **A** across all *c* levels of factor **B**.

$\mu_{.j}$ = population mean for those in the *jth* level of factor **B** across all *r* levels of factor **A**.

		Factor B				Total
		1	2	...	c	
Factor A	1	μ_{11}	μ_{12}		μ_{1c}	$\mu_{1.}$
	2	μ_{21}	μ_{22}		μ_{2c}	$\mu_{2.}$
	...					
	r	μ_{r1}	μ_{r2}		μ_{rc}	$\mu_{r.}$
Total		$\mu_{.1}$	$\mu_{.2}$		$\mu_{.c}$	

Two-Way Analysis of Variance

The global F-test is used to test the null hypothesis that there is no effect of either factor versus the alternative hypothesis that one of the factors has an effect **(there is an effect of either factor A or factor B or both)**.

The test for Factor A is testing the null hypothesis

$H_0 : \mu_{1.} = \mu_{1.} = \dots = \mu_{r.}$ against the alternative that

at least two of the underlying means are different $\mu_{i.} \neq \mu_{j.}$ for some i and j.

The test for Factor B is testing the null hypothesis

$H_0 : \mu_{.1} = \mu_{.1} = \dots = \mu_{.c}$ against the alternative that

at least two of the underlying means are different $\mu_{.i} \neq \mu_{.j}$ for some i and j.

Generally, the procedure for testing follows the procedure for the one-way ANOVA.

Interaction

Besides **the global test and the tests for each factor**, **a test of interaction must be performed** to understand if there is a non-additive effect of the either of the factors on the response.

One must first conduct an interaction test in the two-way ANOVA setting before the **“main effects”** type model can be conducted.

The concept of interactions is very important in statistics and in performing any type of multivariable analysis (including ANCOVA, multivariable regression, and two-way ANOVAs).

Interactions are when the effect of one factor is not constant over levels of another factor.

Interaction

There is **an effect of the drug**, but **no effect of gender**.

- ▷ The presence and absence of an effect can be observed via the use of the **marginal means**.
- ▷ Those on drug B have a higher mean SBP (200) versus those on drug A (100).
- ▷ There does not appear to be an effect of gender (both males and females have an average SBP of 150).

		Drug		Total
		A	B	
Gender	F	100	200	150
	M	100	200	150
Total		100	200	

Interaction

There is **an effect of the drug** and **an effect of gender**.

There also appears to be an effect of gender (males have a higher SBP than females [200 versus 300, respectively]).

Here, the effects of gender are consistent across each drug and the effect of the drug are consistent across levels of gender. We conclude that though males tend to have higher SBPs than women, those on drug A have smaller SBPs than those on drug B. **In this case, there is not an interaction present.**

		Drug		Total
		A	B	
Gender	F	100	300	200
	M	200	400	300
Total		150	350	

Interaction

It appears as if there is **not an effect of drug or gender**, since all of the marginal means are equal.

On closer examination, it appears there is an effect of these factors, **but it varies based on the level of each factor**.

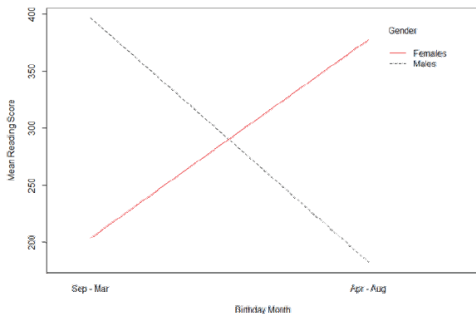
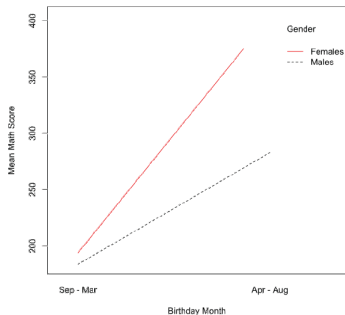
Because the effect of the drug depends on gender (and equivalently, the effect of gender depends on the drug used), **we say that these factors interact**.

		Drug		Total
		A	B	
Gender	F	100	300	200
	M	300	100	200
Total		200	200	

Interaction plots

An educator is interested in the effect of birth month and gender on average reading and math scores for students in the fourth grade. A random sample is taken from each combination.

There is an interaction given that there is not one effect of birth month across both genders.



If there is an interaction between two factors, then the conclusion that we'd get from the **two-way ANOVA are incorrect** at worst and misleading at best.

Before you perform a two-way ANOVA, **you must first check if an interaction is present.**

You do this by including an interaction term in the two-way ANOVA model.

If the **p-value of the interaction term is significant, the two-way ANOVA procedure is not appropriate.** In this case, one **must stratify the analysis by one of the factors and perform separate one-way ANOVAs** for the second factor for each level of the first.

If the interaction p-value is not significant, then one can proceed with the typical two-way ANOVA procedure and assume that the effect of each factor is consistent over the other.

Interaction

Given the importance of detecting interactions and the difference in interpretation of results in the presence of an interaction, alpha levels of 0.10 are often used to check if there is an interaction (as opposed to the more common 0.05 level of significance used in most other settings).

Two-Way Analysis of Variance

```
# Use Anova() function
> Anova([model], type=3)

# First test interaction model
> model = lm(data$response ~ data$group1+ data$group2+ data$group1* data
  $group2)

# Visualize relationship using interaction.plot() function
> interaction.plot(data$group1, data$group2, data$response, col=1:2)

# If p-value for the interaction is not significant, then run regular
  two-way ANOVA
> model = lm(data$response ~ data$group1+ data$group2)
```

An example: Two-Way Analysis of Variance

An exercise physiologist wants to examine whether stretching and wearing ankle weights affect the value of exercise on treadmills. To carry out her study, she recruits subjects who have roughly the same level of physical fitness, and divides them randomly into four groups: with or without ankle weights with or without a stretching period before the exercise.

Using the amount of calories burned as the response, carry out a two-way ANOVA to determine whether stretching and wearing ankle weights have significant effects on exercise.

These are the variables in the data set:

Name	Type	Description
PreStretch	char	stretch group (Stretch, No stretch)
AnkleWeights	char	weights group (Weights, No weights)
Energy	num	calories burned
Speed	num	average speed (in meters per minute)
Oxygen	num	oxygen consumed (in liters)

An example: R commands

```
# Exercise example
exercise <- read.csv("exercise.csv")
exercise
attach(exercise)

#Test interactions
model <- lm(Energy~PreStretch+AnkleWeights+PreStretch*AnkleWeights, data
            =exercise)
summary(model)
Anova(model, type=3)

model1 <- lm(Speed~PreStretch+AnkleWeights+PreStretch*AnkleWeights, data
            =exercise)
summary(model1)
Anova(model1, type=3)

model2 <- lm(Oxygen~PreStretch+AnkleWeights+PreStretch*AnkleWeights,
            data=exercise)
summary(model2)
Anova(model2, type=3)
```

An example: R commands

```
# Generate interaction plots
interaction.plot(PreStretch, AnkleWeights, Energy, col=1:2)
interaction.plot(PreStretch, AnkleWeights, Speed, col=1:2)
interaction.plot(PreStretch, AnkleWeights, Oxygen, col=1:2)

# If interaction is significant, need to stratify (by more of the two
  factors)
stretch <- exercise[which(PreStretch=='Stretch'),]
nostretch <- exercise[which(PreStretch=='No stretch'),]

summary(aov(Energy~AnkleWeights, data=stretch))
summary(aov(Energy~AnkleWeights, data=nostretch))
summary(aov(Speed~AnkleWeights, data=stretch))
summary(aov(Speed~AnkleWeights, data=nostretch))
summary(aov(Oxygen~AnkleWeights, data=stretch))
summary(aov(Oxygen~AnkleWeights, data=nostretch))
```