

## ANOVA and Linear Regression

The purpose of this very brief chapter is to point out the similarities between Analysis of Variance and linear regression. I assume that you are already familiar with the material covered in my chapters on [one-way ANOVA](#) and [simple linear regression](#).

### 1. Review of ANOVA in linear regression

Simple linear regression is concerned with predicting Y from X (when X and Y are both interval or ratio scale variables, and the relationship between them is linear). In linear regression, the total variability of the Y scores is partitioned (or analyzed) into two components:

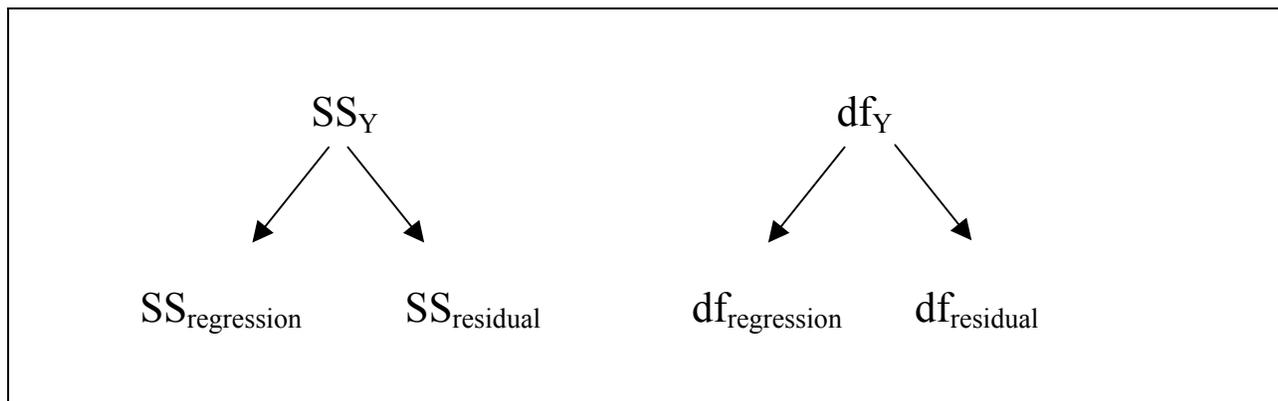
- that which is explained, or accounted for by the linear relationship between X and Y
- that which is left over (residual), or unexplained

We actually partition the “sum of squares Y” ( $SS_Y$ ) rather than the variance of Y. The reason for this, you may recall, is that sums of squares are always additive, whereas variances are not. The portion of  $SS_Y$  that **is** explained by the relationship between X and Y is based on the deviations of predicted scores from the mean of Y:

$$(1) \quad SS_{\hat{Y}} = SS_{\text{regression}} = \sum (\hat{Y} - \bar{Y})^2$$

The portion of  $SS_Y$  that is not explained by the relationship between X and Y (i.e., the **residual** portion) is due to **errors** in prediction, with error defined as actual score minus predicted score:

$$(2) \quad SS_{\text{residual}} = SS_{\text{error}} = \sum (Y - \hat{Y})^2$$



**Figure 1:** Partitioning diagram for ANOVA of simple linear regression.

This partitioning of  $SS_Y$  (and  $df_Y$ ) is illustrated in Figure 1. The F-test in linear regression is calculated as follows:

$$(3) \quad F_{(p, N-p-1)} = \frac{MS_{regression}}{MS_{residual}} = \frac{(SS_{regression} \div df_{regression})}{(SS_{residual} \div df_{residual})}$$

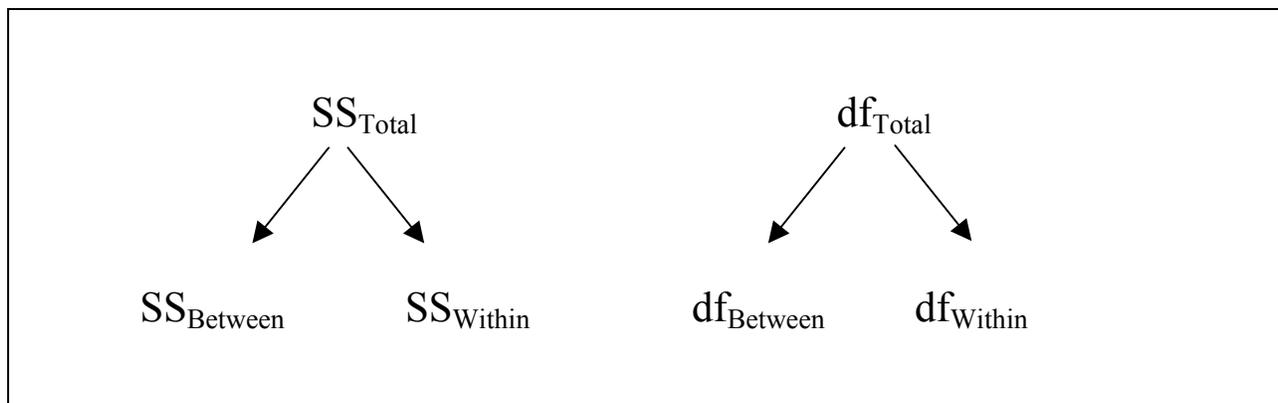
In the case of simple linear regression (i.e., only 1 predictor variable),  $df_{regression} = 1$ , and  $df_{residual} = N-2$  (where  $N$  = the total number of paired  $X$ - $Y$  scores). In general,  $df_{regression} = p$ , where  $p$  = the number of predictor variables; and  $df_{residual} = N-p-1$ .

Finally, you may recall that the **coefficient of determination**,  $r^2$ , represents the proportion of variability in  $Y$  that is accounted for by the linear relationship between  $X$  and  $Y$ :

$$(4) \quad r^2 = \frac{SS_{regression}}{SS_Y} = 1 - \frac{SS_{residual}}{SS_Y}$$

## 2. One-way ANOVA: Same concept, slightly different terminology

Exactly the same kind of partitioning of the overall variance occurs in one-way ANOVA. The only difference is that it occurs in the context of  $k$  independent groups of subjects ( $k \geq 2$ ). Many textbook authors use  $\mathbf{X}$  rather than  $\mathbf{Y}$  to represent the dependent variable for ANOVA problems. For consistency, however, I will stick with  $\mathbf{Y}$ .



**Figure 2:** Partitioning diagram for one-way ANOVA.

The partitioning of  $SS_{Total}$  (and  $df_{Total}$ ) in one-way ANOVA is illustrated in Figure 2. The F-test in one-way ANOVA is calculated as follows:

$$(5) \quad F_{(k-1, N-k)} = \frac{MS_{Between}}{MS_{Within}} = \frac{(SS_{Between} \div df_{Between})}{(SS_{Within} \div df_{Within})}$$

In one-way ANOVA,  $df_{Between} = k-1$ , where  $k$  = the number of groups; and  $df_{Within} = N-k$  (where  $N$  = the total number of scores).

Finally, one can calculate a measure that is analogous to  $r^2$  in simple linear regression. It is called *eta-squared*, or the *correlation ratio*. The Greek letter *eta* looks like this:  $\eta$ . So the symbol for eta-squared is  $\eta^2$ . It is computed as shown below, and is equal to the proportion of the total variance of the Y-scores that is explained by the independent variable.

$$(6) \quad \eta^2 = \frac{SS_{Between}}{SS_{Total}} = 1 - \frac{SS_{Within}}{SS_{Total}}$$

The correspondences between ANOVA of regression and one-way ANOVA are summarized in Table 1, which due to its width is shown in the appendix.

### 3. Regression and ANOVA: Special cases of the General Linear Model

Given the similarities between regression and ANOVA, it may not surprise you to learn that any so-called ANOVA problem can be analyzed using linear regression. Section 2.3 of my chapter on [ANCOVA](#) provides details on how this works. The reason for this is that linear regression and ANOVA are both just special cases of the *general linear model* (or GLM for short).

All linear models are of the same basic form. That is, the predicted, or modeled value of some outcome variable Y is equal to a linear combination of predictor variables plus an error term. *Linear combination* just means that each of the predictor variables is multiplied by a coefficient. In symbols:

$$(7) \quad \hat{Y} = c_1X_1 + c_2X_2 + \dots + error$$

I'm sure you're all familiar with the expression, *Genius is 1% inspiration and 99% perspiration*. That expression is an example of a linear combination. If we let Y = Genius,  $X_1$  = inspiration, and  $X_2$  = perspiration, we can write the equation as follows:

$$(8) \quad \hat{Y} = 0.01(X_1) + 0.99(X_2) + error$$

If all linear models have this same basic form, why do we have different techniques such as linear regression, ANOVA, and ANCOVA. The reason is that linear models do vary in terms of whether the predictor variables are categorical or interval/ratio scale, or some combination of the

two.<sup>1</sup> As it turns out, for particular combinations of predictor variable types, the calculations can be reduced to a relatively simple form—a shortcut, if you will. Bear in mind that when Karl Pearson and Ronald F. Fisher (and others) were doing the pioneering work on these techniques, folks did not have the computing power available to them that we enjoy today. So simplifications and shortcuts for specific types of models were very useful. I suspect that if the computing power had been available back then, we would probably not talk about regression, ANOVA, and ANCOVA as distinct techniques—we would probably just describe them all as particular types of linear models. The characteristics of some of the more common linear models you may encounter are summarized in Table 2.

**Table 2:** Characteristics of some common types of linear models.

<b>Model Name</b>	<b>Dependent Variable(s)</b>	<b>Independent Variable(s)</b>
Simple linear regression	<ul style="list-style-type: none"> <li>• One DV with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• One IV with interval/ratio scale properties</li> </ul>
Multiple regression	<ul style="list-style-type: none"> <li>• One DV with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• Two or more IVs, all of which have interval/ratio scale properties</li> </ul>
One-way ANOVA	<ul style="list-style-type: none"> <li>• One DV with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• One categorical IV (i.e., a grouping variable)</li> </ul>
Two-way ANOVA	<ul style="list-style-type: none"> <li>• One DV with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• Two categorical IVs</li> </ul>
One-way MANOVA	<ul style="list-style-type: none"> <li>• Two or more DVs, all with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• One categorical IV</li> </ul>
One-way ANCOVA	<ul style="list-style-type: none"> <li>• One DV with interval/ratio scale properties</li> </ul>	<ul style="list-style-type: none"> <li>• One categorical IV</li> <li>• One interval/ratio IV (the <i>covariate</i>)</li> </ul>

---

<sup>1</sup> The type of outcome variable can also vary (e.g., it can be dichotomous), but for now, I am restricting myself to cases where the outcome variable (Y) has interval/ratio scale properties.

## APPENDIX

Table 1: Correspondences between ANOVA of regression and one-way ANOVA.

Symbol in ANOVA of Linear Regression	Description	Corresponding symbol(s) for One-way ANOVA	Description
$\bar{Y}$	<ul style="list-style-type: none"> <li>The Grand Mean</li> </ul>	$\bar{Y}_{Grand}, \bar{Y}_{Total}, \bar{Y}.$	<ul style="list-style-type: none"> <li>The Grand Mean</li> </ul>
$\hat{Y}_i$ or $Y'_i$	<ul style="list-style-type: none"> <li>The predicted score for the <math>i^{th}</math> person</li> </ul>	$\bar{Y}_j$	<ul style="list-style-type: none"> <li>Mean of the <math>j^{th}</math> sample</li> <li>the predicted score for a person in the <math>j^{th}</math> sample</li> </ul>
$SS_Y$ or $SS_{Total}$	<ul style="list-style-type: none"> <li>Sum of the squared deviations about the grand mean</li> </ul>	$SS_{Total}$	<ul style="list-style-type: none"> <li>Sum of the squared deviations about the grand mean</li> </ul>
$SS_{regression}$	<ul style="list-style-type: none"> <li>Sum of squared deviations of predicted scores about the grand mean</li> </ul>	$SS_{Between-groups}$ or $SS_{Treatment}$	<ul style="list-style-type: none"> <li>Sum of squared deviations of group means about the grand mean</li> </ul>
$SS_{residual}$	<ul style="list-style-type: none"> <li>Sum of squared deviations of actual scores around predicted scores</li> </ul>	$SS_{Within-groups}$ or $SS_{error}$	<ul style="list-style-type: none"> <li>Sum of squared deviations of actual scores around group means</li> </ul>