

We're now able to fit a MLR model (i.e., estimate the β 's), we can interpret the relationships (via the β 's), we can test hypotheses about the β 's and we can calculate CI's and PI's for expected/predicted response values.

We should probably decide whether the model fits the data well... We'll do this with an Analysis of Variance (ANOVA).

Analysis of Variance (ANOVA)

Variability in the response variable (the y values) can be quantified by how far they deviate from their mean. Measures of variation in the y 's typically contain the term

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

which we define as the total sum of squares

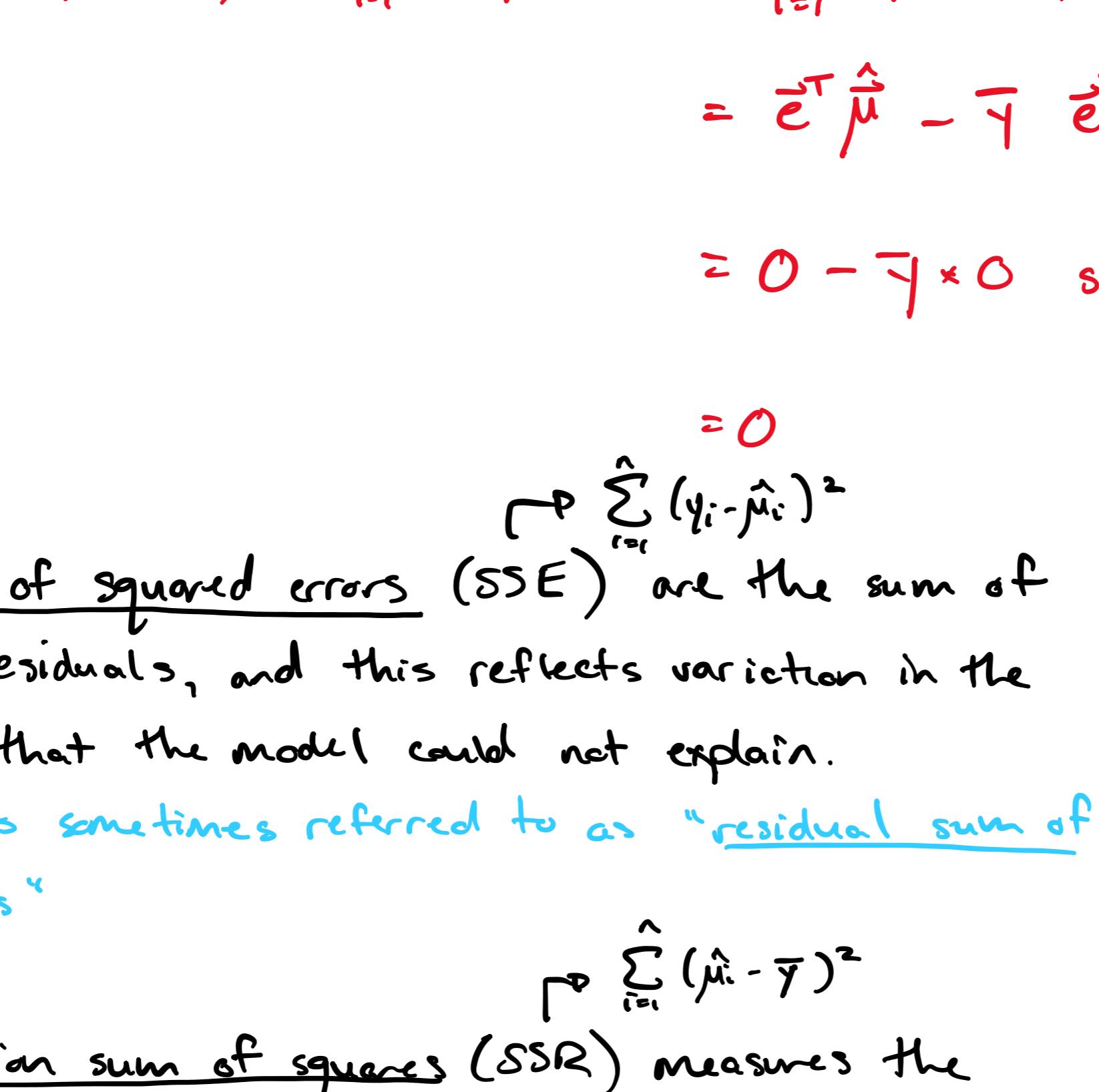
* Notice $S_y = \sqrt{\frac{SST}{n-1}}$

The objective of ANOVA is to partition the total variation in y (SST) into two parts:

- (1) Variation accounted for by the model
- (2) Variation not accounted for by the model

Notice that $y_i - \bar{y} = (y_i - \hat{\mu}_i) + (\hat{\mu}_i - \bar{y})$, $i=1, 2, \dots, n$

Although this is true no matter how many explanatory variables we have, we'll visualize it in SLR case.



Let's now look at this partition in terms of the sums of squares:

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n [(y_i - \hat{\mu}_i) + (\hat{\mu}_i - \bar{y})]^2 \\ &= \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 + \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{\mu}_i)(\hat{\mu}_i - \bar{y}) \end{aligned}$$

$$= \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 + \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2$$

$$= SSE + SSR$$

$$\sum_{i=1}^n (y_i - \hat{\mu}_i)(\hat{\mu}_i - \bar{y}) = \sum_{i=1}^n e_i (\hat{\mu}_i - \bar{y}) = \sum_{i=1}^n e_i \hat{\mu}_i - \bar{y} \sum_{i=1}^n e_i$$

$$= \vec{e}^T \hat{\mu} - \bar{y} \vec{e}^T \vec{1}$$

$$= 0 - \bar{y} \times 0 \text{ since } \vec{e} \text{ is orthogonal to both } \hat{\mu} \text{ and } \vec{1}$$

$$= 0$$

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{\mu}_i)^2$$

The sum of squared errors (SSE) are the sum of squared residuals, and this reflects variation in the response that the model could not explain.

* SSE is sometimes referred to as "residual sum of squares"

$$\sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2$$

The regression sum of squares (SSR) measures the amount of variation in the response that the model is able to explain.

This variation decomposition is typically displayed in an ANOVA table.

Source	Sum of Squares	df	Mean Squares	F
Regression	SSR	p	$MSE_R = \frac{SSR}{p}$	$F_0 = \frac{MSE_R}{MSE_E}$
Error	SSE	n-p-1	$MSE_E = \frac{SSE}{n-p-1}$	
Total	SST	n-1		

$$* df(SSR) + df(SSE) = df(SST)$$

$$\hat{\sigma}^2$$

* the table also contains "mean squares" which are the sums of squares divided by their respective degrees of freedom

test statistic associated with the test overall significance in the linear regression (will come back to this)

To answer the question "Does the model fit the data well?" we can use the coefficient of determination (R^2). This metric quantifies the proportion of response variation accounted for by the model.

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \in [0, 1]$$

$$\uparrow \text{because } SST = SSR + SSE$$

Intuitively larger values are better than small values.

The more variation the model explains the larger SSR will be (and hence the smaller SSE will be). What this means practically is that the fitted values lie close to true values and hence our model fits the data well.

PGA Example

Source	df	Sum Sq.	Mean Sq.	F
Regression	1	2049.0	2049.0	115.31
Error	194	3447.3	17.77	
Total	195	5496.3		

$$R^2 = \frac{2049.0}{5496.3} = 0.3728$$

$$5496.3$$

Therefore, the model (driving distance) explains 37.28% of the variation in the response (driving accuracy).