2/9/2023

**Multiple linear regression**
Multiple linear regression is an extension of simple linear regression, but here we are using more than one independent variable to predict one output variable.  Typically, our model equation takes the form:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_k x_k + \epsilon$$

Like the simple linear model, the $\beta$'s represent the true regression coefficients of the population which we estimate with $b$'s. The variables in the model are all linear, and the error term at the end has a mean of zero and a constant standard deviation, which we can estimate from the residuals.  The model in multiple dimensions is a (hyper)plane – when using two independent variables, it is a regular plane which we can graph and examine explicitly. When there are more variables, we can no longer look at the graphs, but we can still make predictions and do our other types of analysis. The linear equation here is sometimes referred to as a first-order model.

When we have more than one slope variable the ANOVA analysis for the model is no longer equivalent to the slope test since we have to look at more than one slope. We should now interpret the ANOVA, full-model analysis, as testing whether any coefficient for the variables in the model are non-zero, similar to a more traditional ANOVA. If the P-value is too high (greater than the significance level), then we conclude that none of the coefficients are non-zero.  If the P-value is below the significance level, then we can conclude that at least one coefficient is non-zero.

From that point, we conduct tests on the individual coefficients.  Testing information (test-statistic, standard error, P-value) is included in the summary results in R.  We must then consider which if the variables is statistically significant. It may be that they all are, or only some of them.  Later we will develop model selection strategies for building models with coefficients that are all statistically significant. For our initial discussion, we are going to consider the steps for building our initial model, testing parameters, and other new things we need to consider in the multi-variable case.

Each of our multiple-variable models may have different numbers of variables, and each case will require slightly different summation equations.  However, the equation will use for the linear algebra approach does not change.  We are, therefore, going to use that method with a small example to illustrate how the setup changes as the number of variables increases. We'll use technology to solve the models for us in any real world context.

Let's consider a situation where we have two variables that are independent which are predicting a third variable. Our data will then come in ordered triples, and our linear model will have the form

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

Let's suppose our data is $\{(1,2,10), (2,4,15), (2,6,18), (3,7,24), (4,7,27)\}$.  The first coordinate is $x_1$, the second coordinate is $x_2$, and the third coordinate is $y$.  Replace these into our model equation.

$$\beta_0 + \beta_1(1) + \beta_2(2) = 10$$
$$\beta_0 + \beta_1(2) + \beta_2(4) = 15$$
$$\beta_0 + \beta_1(2) + \beta_2(6) = 18$$
$$\beta_0 + \beta_1(3) + \beta_2(7) = 24$$

$$\beta_0 + \beta_1(4) + \beta_2(7) = 27$$

The coefficients of our $\beta$'s go into our A matrix, the $\beta$'s go into the B matrix, and the constants go into the Y matrix.

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 4 \\ 1 & 2 & 6 \\ 1 & 3 & 7 \\ 1 & 4 & 7 \end{bmatrix}, B = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}, Y = \begin{bmatrix} 10 \\ 15 \\ 18 \\ 24 \\ 27 \end{bmatrix}$$

We solve this using the normal equation

$$A^T A B = A^T Y$$

Or the solution matrix

$$B = (A^T A)^{-1} A^T Y$$

We find that $B \approx \begin{bmatrix} 3.176 \\ 3.706 \\ 1.294 \end{bmatrix}$, which gives us the equation $\hat{y} = 3.176 + 3.706 x_1 + 1.294 x_2$.

Let's think about interpretation a moment.  The constant is the predicted value of $y$ when both $x_1$ and $x_2$ are equal to zero. As with the simple linear model, this may not be possible. If zero is far outside the domain of even one variable, then the intercept may be meaningless.  So, interpret with care.  The coefficient of each variable can be treated as regular slopes if all other variables in the equation are held constant.  For instance, the value of $y$ increases by 3.706 for each one unit increase in $x_1$ if $x_2$ remains the same.  Likewise, the $y$ value increases by 1.294 for each one unit increase in $x_2$ if $x_1$ remains constant.

One way we analyzed the relationship between variables in the simple linear model was using correlation. But correlation on works on pairs of variables, not three or more.  There is another way to analyze models with multiple variables, using $R^2$. In other contexts, it's referred to as the coefficient of determination.  In multiple variable models, it may be referred to as the coefficient of multiple determination.  It really is the same thing.  In the simple linear case, we can square the correlation to obtain this value, but in this case, we need another method to arrive at this number.

$$R^2 = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

The numerator is the difference between predicted values and the mean, squared. The denominator is the difference between the original observations and the mean, squared.  The denominator can be thought of as the variability in the original observed values, and the numerator in the variability in the predicted values.  The ratio here can be interpreted as the amount of the variability in the $y$ value that can be accounted for by the model relationship.  That means that a high $R^2$ means that the model is a good fit for the data because the relationship account for most of the changes in y-values.  If there is a low $R^2$ values, then the model accounts for very little of the variability and does little to improve our

predictions over just using the mean.  Because the value must be between 0 and 1, it is often expressed as a percentage.  When we do a model summary in R, $R^2$ is one of the reported values.

The value $1 - R^2$ is sometimes thought of as the fraction of the original variability left in the residuals.

Our example model included only two variables, but we can include many more. How do we know if we've included too many? One way is to adjust the $R^2$ value for the number of variables included in the model. In statistics, we generally want to follow the principle of parsimony, which says that we want the simplest model possible that produces the best predictions.  We can overfit models if we use too many variables and so the adjusted $R^2$ is a way of considering this.  $adj\text{-}R^2 =$

$$R_a^2 = 1 - \frac{n-1}{n-(k+1)}\left(\frac{\sum(y_i - \hat{y}_i)^2}{\sum(y_i - \bar{y})^2}\right)$$

Here, $k$ is the number of variables used in the model, and $n$ is the sample size. The adjusted $R^2$ values will generally be smaller than $R^2$, but as we add variables, there will come a point when adding variables stops adding much to $R^2$, and may even result in decreasing the adjusted $R^2$. This is a sign that adding the extra variable is not adding enough new information or power to the model to justify using the extra variables.

The adjusted $R^2$ is also often included in our model summaries when we do regression in R.

We can use the $R^2$ to obtain what is called $R$ the multiple correlation coefficient. It is the positive square root of the coefficient of determination.  You can use the same range of values we used for correlation to assess the strength of the model.

As with the simple linear model, we can construct confidence intervals, and prediction intervals. We can conduct hypothesis testing on each variable in a model, and construct confidence intervals on each coefficient.  These procedures follow the same general methods we used in the simple linear regression model.

One additional concern that we have with multiple variable models that we did not have before is that our "independent" variables may not be truly independent.  We would prefer that our independent variables not be highly correlated with each other. This issue is sometimes referred to as collinearity. This will be one of the tests we'll want to conduct when we assess our models.

Let's look at a model using our mtcars dataset.  Let's look at a model of mpg using the disp(lacement) variable and the wt(weight) variable. The summary output looks like this.

Call:
lm(formula = mpg ~ disp + wt, data = mtcars)

Residuals:

| Min | 1Q | Median | 3Q | Max |
| --- | --- | --- | --- | --- |
| -3.4087 | -2.3243 | -0.7683 | 1.7721 | 6.3484 |

Coefficients:

|  | Estimate | Std. Error | t value | Pr(>\|t\|) |
|---|---|---|---|---|
| (Intercept) | 34.96055 | 2.16454 | 16.151 | 4.91e-16 *** |
| disp | -0.01773 | 0.00919 | -1.929 | 0.06362 . |
| wt | -3.35082 | 1.16413 | -2.878 | 0.00743 ** |

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

Residual standard error: 2.917 on 29 degrees of freedom
Multiple R-squared: 0.7809,     Adjusted R-squared: 0.7658
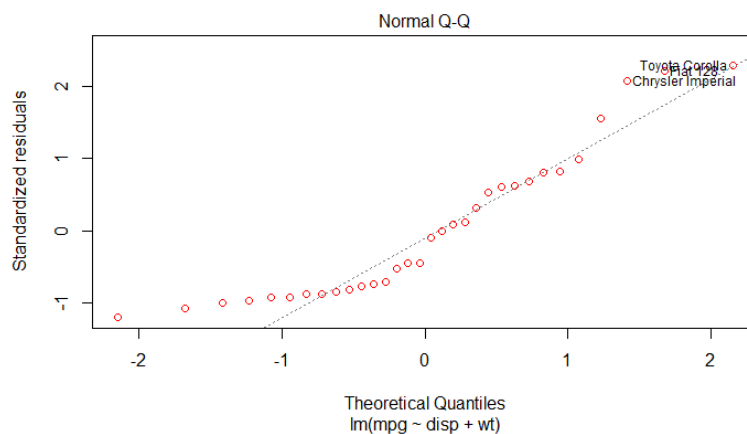F-statistic: 51.69 on 2 and 29 DF,  p-value: 2.744e-10

Let's first consider our full model.  Look at the last line of the summary output.  This is the result of the ANOVA test, and we have a P-value that is very small. This indicates that at least one of the variable coefficients in the model is non-zero. So something in this model helps to improve our predictions of mpg.
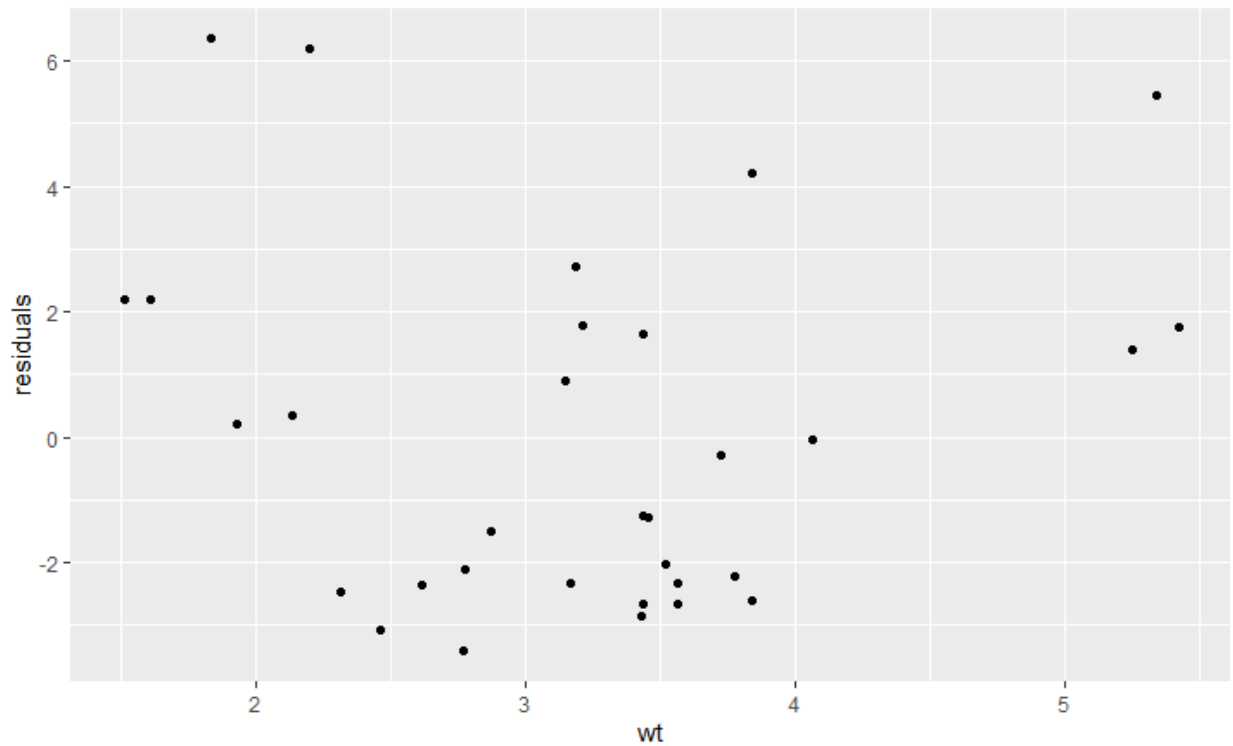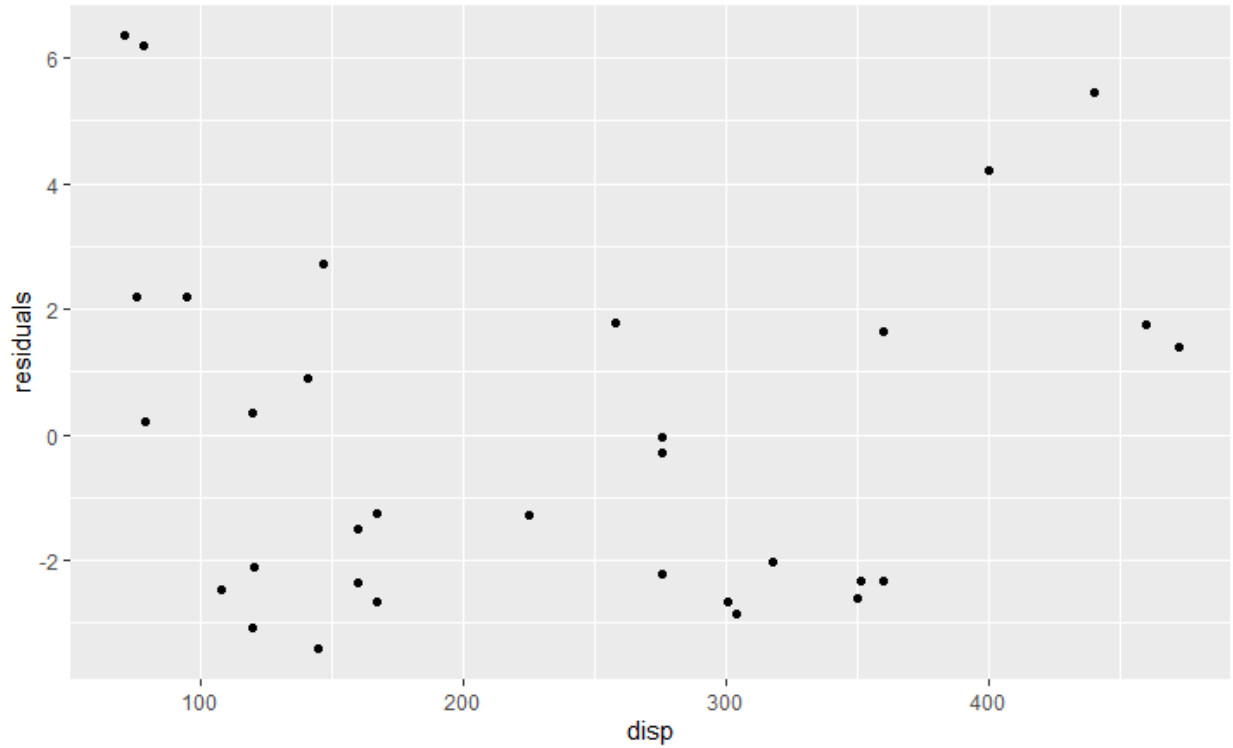
The multiple $R^2$ is about 78% (we've reduced the variability by this much). The adjusted-$R^2$ is similar but a bit smaller.

Let's look at our table of coefficients.  The P-value for the intercept is quite small, so it's not zero.  The P-value for the disp variable is 0.06, which is significant at the 10% level, but too high at the 5% level.  The coefficient for weight, however, the P-value is 0.007 which is much less than 5%, so this variable should be kept in the model. If we think that the disp variable lacks significance, we can remove it from the model and rerun the test.

We should check for collinearity in our "independent" variables. We can do this by testing the correlation between them. We find that the Pearson correlation is 88%, which is probably why both variables are not statistically significant. Only one of them is contributing new information.  We may wish to test the simple linear models with each variable to see which does a better job predicting mpg on their own.

We should look at the residual plots and the normality plot for the residuals. When we have more than one x-variable, we plot the residuals against each variable separately.



Normal Q-Q
lm(mpg ~ disp + wt)

The normality plot indicates that the residuals are not very normal. The weight residual plot looks a bit better than the disp residual plot. This one seems to be much more widely dispersed on the ends and more positive, while in the middle the values are smaller and more likely to be negative. An effect that is not as extreme in the second plot. These may be signs of lack of linearity or a lack of constant variance. It suggests that our model assumptions have not been met.

It may be that we are willing to trade off some of these assumptions for a easy to understand model that is better than nothing. However, we can return to this when we have developed more tools for dealing with nonlinear models.

In the next lecture we will look more closely at outliers and influential points: detecting them, and what to do with them.

References:

1. https://faculty.ksu.edu.sa/sites/default/files/probability_and_statistics_for_engineering_and_the_sciences.pdf
2. http://www.sthda.com/english/wiki/correlation-test-between-two-variables-in-r