

# Assumptions Underlying the Regression Model

2019-09-24

## Introduction and Research Question

In this set of notes, we will return again to examine the question of whether time spent on homework is related to GPA using the *keith-gpa.csv* data (see the [data codebook](#)). To begin, we will load several libraries and import the data into an object called *keith*.

## Preparation

```
# Load libraries
library(broom)
library(corr)
library(dplyr)
library(educate)
library(ggplot2)
library(readr)
library(tidyr)

keith = read_csv(file = "~/Documents/github/epsy-8251/data/keith-gpa.csv")
head(keith)
```

```
# A tibble: 6 x 3
  gpa homework parent_ed
<dbl>   <dbl>   <dbl>
1    78         2       13
2    79         6       14
3    79         1       13
4    89         5       13
5    82         3       16
6    77         4       13
```

```
# Fit simple regression model
lm.1 = lm(gpa ~ 1 + homework, data = keith)
```

## Model Assumptions

There are five primary assumptions we make about the regression model in order for the results we obtain from fitting this model (e.g., coefficient estimates,  $p$ -values, CIs) to be valid.

- Linearity
- Independence
- Normality
- Homoskedasticity (homogeneity of variance)

We will be more specific about each of these, but it can be instructive to examine a visual representation of some of these assumptions. Imagine that we had the population of  $(x, y)$  values and we plotted them and fitted a regression to them. That picture would look like this,

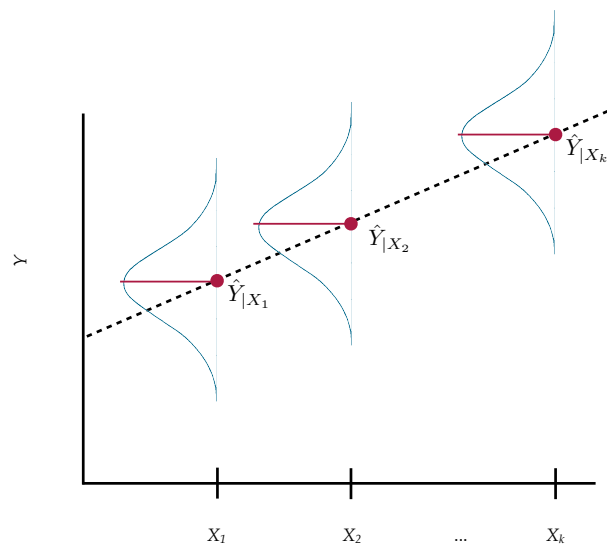


Figure 1. A visual depiction of the simple regression model's assumptions.

In Figure 1, the normal distributions being shown are the distribution of  $Y$ -values at each value of  $X$ ; the distribution of  $Y$  is conditioned on  $X$ , and is thus called a conditional distribution. Although only three of these distributions are shown in the figure, there is a conditional distribution for EVERY value of  $X$ . Now that we understand this picture, we can expand upon the regression assumptions listed earlier.

- **Linearity:** The linearity assumption implies that the MEAN values of  $Y$  from all the conditional distributions all fall on the same line. If this is the case, we would say that the conditional mean  $Y$ -values are linear.
- **Independence:** This is not shown in the figure. The assumption is that each  $Y$ -value in a particular conditional distribution is independent from every other  $Y$ -value in that same distribution.
- **Normality:** This assumption indicates that every one of the conditional distributions of  $Y$ -values is normally distributed.
- **Homoskedasticity:** This is the homogeneity of variance assumption. It says that the variance (or standard deviation) of all of the conditional distributions is exactly the same.

## Assumptions are Really About the Residuals

Note that so far we have stated these assumptions in terms of the  $Y$ -values and the conditional distributions of  $Y$ . Technically, all model assumptions (for regression, ANOVA,  $t$ -test, etc.) all refer to the residuals. Think about how we compute the residuals:

$$\epsilon_i = Y_i - \hat{Y}_i$$

In Figure 1, the  $\hat{Y}_i$  value is the  $Y$ -value that corresponds to the point on the line. If we transform every  $Y$ -value in the population, from Figure 1, to an  $\epsilon$  value, and re-plot them, the visual depiction now looks like this.

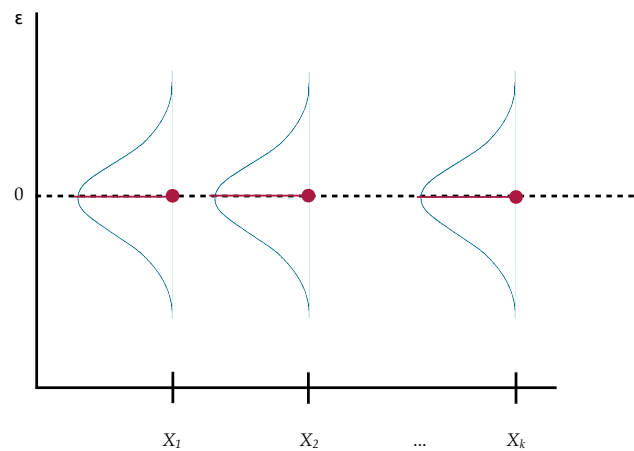


Figure 2. A visual depiction of the simple regression model's assumptions about the residuals.

So if we restate the assumptions in terms of the residuals and the conditional distributions of the residuals,

- **Linearity:** The MEAN value of each of the conditional distributions of the residuals is zero.
- **Independence:** Again, this is not shown in the figure. The assumption is that each residual value in a particular conditional distribution is independent from every other residual value in that same distribution.
- **Normality:** This assumption indicates that each of the conditional distributions of residuals is normally distributed.
- **Homoskedasticity:** The variance (or standard deviation) of all of the conditional distributions of residuals is exactly the same.

These assumptions can also be expressed mathematically as,

$$\epsilon_{i|X} \sim \text{i.i.d } \mathcal{N}(0, \sigma^2)$$

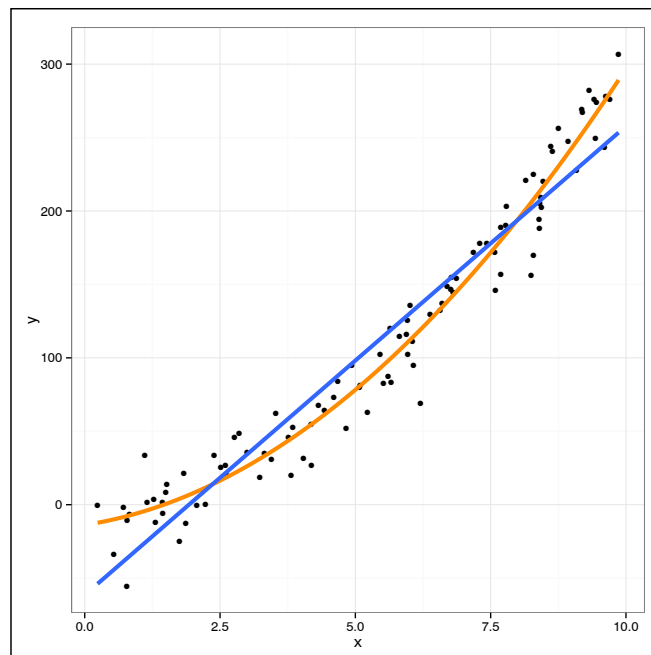
The “i.i.d” stands for *independent and identically distributed*. The mathematical expression says the residuals conditioned on  $X$  are independent and identically normally distributed with a mean of 0 and some variance, represented as  $\sigma^2$ .

## Evaluating the Regression Model's Assumptions

Before beginning to evaluate the assumptions using our data, it is important to point out that the assumptions are about the population's residuals. Because in most analyses, we only have a sample of data, we can never really evaluate these assumptions. We can only offer a guess as to whether they are tenable given the data we see. The strongest arguments for justification for meeting any of the model's assumptions is a theoretical argument based on existing literature in the discipline.

### Linearity

The linearity assumption is critical in specifying the structural part of the model. Fitting a linear model when the TRUE relationship between  $X$  and  $Y$  is non-linear may be quite problematic. Coefficients may be wrong. Predictions may also be wrong, especially at the extreme values for  $X$ . More importantly, mis-specified models lead to misinformed understandings of the world.



*Figure 3.* The plot shows the differences between the true non-linear model (orange) and a mis-specified linear fitted model (blue). Using the linear fitted model to make predictions would be quite misleading, especially at extreme values of  $X$ .

Notice that when a linear model is fitted to data generated from a non-linear function (as in Figure 3) that the data are consistently above, or below the line, depending on the  $X$ -value. This type of pattern would be evidence that the linearity assumption is not tenable. When evaluating this assumption, we want to see data in the scatterplot that is “equally” above and below the line at each value of  $X$ .

### Independence

The definition of independence relies on formal mathematics. Loosely speaking a set of observations is independent if knowing that one observation is above (or below) the mean value in a conditional distribution conveys no information about whether any other observation in the same distribution is above (or below) its mean value. If observations are not independent, we say they are dependent or correlated.

Using random chance in the design of the study, to either select observations (random sampling) or assign them to levels of the predictor (random assignment) will guarantee independence of the observations. Outside of this, independence is often difficult to guarantee, and often has to be a logical argument.

There are a few times that we can ascertain that the independence assumption would be violated. These instances often result from aspects of the data collection process. One such instance common to social science research is when the observations (i.e., cases, subjects) are collected within a physical or spatial proximity to one another. For example, this is typically the case when a researcher gathers a convenience sample based on location, such as sampling students from the same school. Another violation of independence occurs when observations are collected over time (longitudinally), especially when the observations are repeated measures from the same subjects.

One last violation of independence occurs when the observation level used to assign cases to the different predictor values (e.g., treatment or control) does not correspond to the observation level used in the analysis. For example, in educational studies whole classrooms are often assigned to treatment or control. That means that the cases used in the analysis, in order to satisfy the independence assumption, would need to be at the classroom level (e.g., the cases would need to be classroom means), not individual students. This can be deleterious for sample size.

If the independence assumption is violated, almost every value you get in the `tidy()` and `glance()` output—the standard errors, *t*-values, *p*-values, *F*-statistics, residual standard errors—are wrong. If you suspect that you have violated the independence assumption, then you will need to use a method (not OLS regression) that accommodates non-independence. (We cover some of these methods in EPsy 8252.)

## Normality and Homoskedasticity

The assumptions about normality and homoskedasticity are about the distribution of errors at each level of *X*. Both of these assumptions are less critical than the assumptions of linearity and independence. It is only problematic for the OLS regression results if there are egregious violations of the distributional assumptions. In general, if these assumptions are only minorly violated, the results of the OLS regression are still valid; we would say the results from an OLS regression are *robust* to violations of normality and homoskedasticity. Even if the violations are bad, there are many transformations of the data that can alleviate those problems. We will cover some of those transformations later in the course.

## Empirically Evaluating the Assumptions

We can use the data to empirically evaluate the assumptions of linearity, normality, and homoskedasticity. (The assumption of independence is difficult to evaluate using the data, and is better left to a logical argument.) Recall that the assumptions are about the residuals. To compute the residuals, we will use the `augment()` function from the **broom** package. We will also write those results into an object so we can compute on it later.

```
# Augment the model to get residuals
out_1 = augment(lm.1)

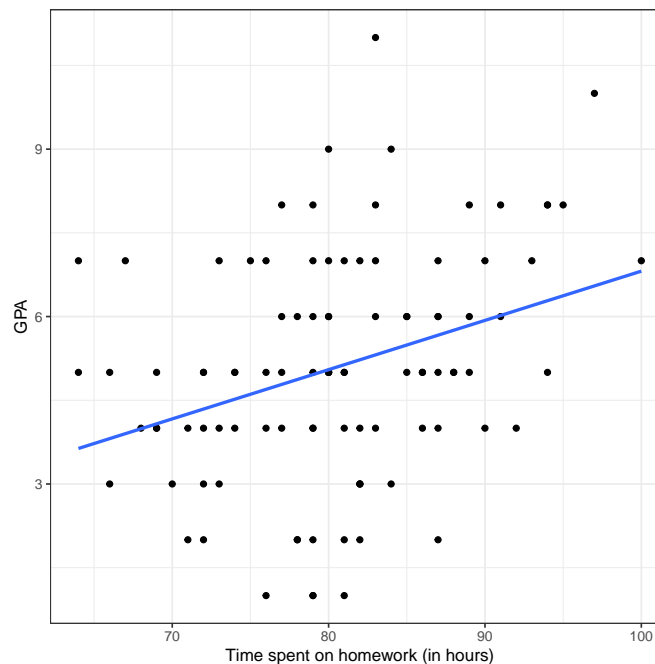
# View augmented data
head(out_1)
```

```
# A tibble: 6 x 9
  gpa homework .fitted .se.fit .resid .hat .sigma .cooksd .std.resid
  <dbl>   <dbl>   <dbl>   <dbl> <dbl> <dbl> <dbl>   <dbl>   <dbl>
1    78     2    76.7    1.31   1.28 0.0328  7.28 0.000550    0.180
2    79     6    81.6    0.792 -2.57 0.0120  7.27 0.000776   -0.358
3    79     1    75.5    1.62   3.50 0.0500  7.27 0.00646    0.495
4    89     5    80.4    0.725   8.64 0.0100  7.22 0.00728    1.20
5    82     3    77.9    1.04   4.07 0.0204  7.26 0.00336    0.568
6    77     4    79.1    0.820  -2.15 0.0128  7.27 0.000579   -0.298
```

The assumptions are also really about the population of residuals. Because we do not have the entire population of residuals (we obtained our residuals by fitting a regression to a sample of data) we do not expect that our residuals will actually meet the assumptions perfectly (remember, sampling error). Examining the sample residuals, is however, a reasonable way to evaluate the tenability of assumptions in practice. We just have to keep in mind that the sample residuals may deviate a bit from these assumptions.

## Evaluating the Tenability of Linearity

The linearity assumption can initially be evaluated by examining the scatterplots of the outcome versus each predictor. If the relationship is linear we would expect to see roughly half of the observations above the fitted line and half the observations below the fitted line at each value of  $X$ .



*Figure 4.* Scatterplot of GPA versus time spent on homework for 100 students. The plot indicates that linearity may be a tenable based on the pattern observed in these data.

Based on the plot, the linearity assumption seems tenable for these data. However, we will double-check that this holds by also examining a plot of the residuals versus the predictor ( $X$ ). We can use the augmented data to create this plot. In this plot, we want to see that the residuals at each value of  $X$  are equally distributed above and below the  $Y = 0$  line. There should also be no visible pattern in the scatterplot, since the correlation of the predictor with the residuals should be 0.

```
ggplot(data = out_1, aes(x = homework, y = .resid)) +  
  geom_point() +  
  theme_bw() +  
  geom_hline(yintercept = 0) +  
  xlab("Time spent on homework (in hours)") +  
  ylab("Residuals")
```

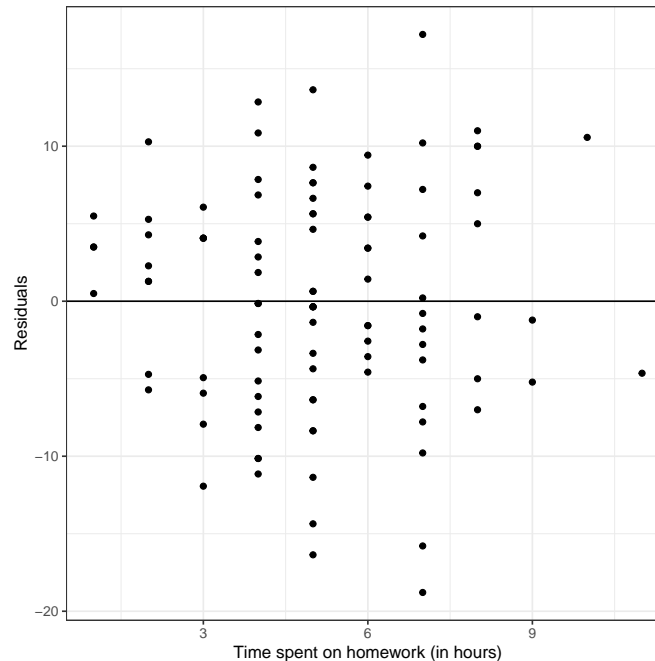


Figure 5. Scatterplot of the residuals versus education level. A horizontal line at  $y = 0$  has been added to the plot to indicate the expected residual value of 0.

Based on the plot of the residuals versus the predicted values the linearity assumption seems tenable. We are seeing that about half the residuals are above the  $Y = 0$  line and half are below the line at each value of  $X$ .

### Evaluating the Tenability of Homoskedasticity

To evaluate homoskedasticity, we want to look at the same plot of the residuals versus predictor values to see if the variation in residuals seems consistent at different values of  $X$  (time spent on homework). In other words, does the range of the residuals look about the same across the different education values?

In general, it seems that the homoskedasticity assumption is also tenable. The range of the residuals looks pretty consistent across the different conditional distributions. Although some ranges look larger (e.g., homework = 7) and some look smaller (e.g., homework = 9) this is probably just a function of sampling variation and not a systematic difference in variation.

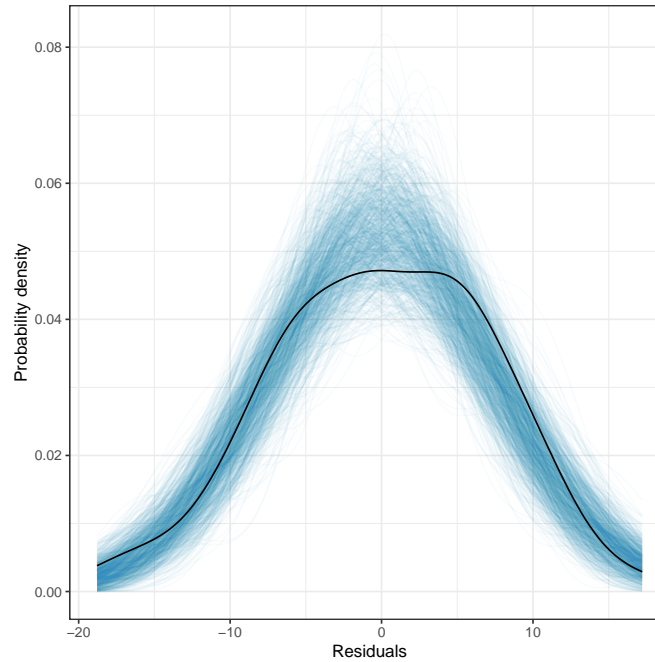
### Evaluating the Tenability of Conditional Normality

The conditional normality assumption requires that the residuals in the population are normally distributed at each value of  $X$ . This can be difficult to evaluate in a scatterplot. Perhaps we could look at a density plot of the residuals at each  $X$  value, but doing this has its own set of problems. For example, there are only two residuals at  $X = 9$ . It would be difficult to evaluate whether they come from a population that was normally distributed or not. This issue of data scarcity is a problem at many of the  $X$  values.

In practice, researchers combine all the residuals together into a single distribution (the *marginal distribution*) and examine whether this distribution is roughly normal, rather than examining the separate conditional distributions. This is easier to do, but is not satisfactory. Even if the marginal distribution seems normal it does not guarantee that the conditional distributions are all normally distributed. Luckily, recall that the normality assumption is not quite as critical, so in general, examining the marginal distribution is a reasonable alternative to not being able to actually examine the conditional distributions.

We will use the `stat_watercolor_density()` function from the `educate` package with the argument `model="normal"` to plot a confidence envelope of where we would expect to see the density IF the population were normally distributed. We can then overlay this with a density plot of the actual residuals using `stat_density()`.

```
ggplot(data = out_1, aes(x = .resid)) +  
  stat_watercolor_density(model = "normal") +  
  stat_density(geom = "line") +  
  theme_bw() +  
  xlab("Residuals") +  
  ylab("Probability density")
```



*Figure 6.* A density plot of the residuals from the fitted regression model. The bootstrapped confidence envelope (in blue) shows the sampling variation in density expected under the normality assumption.

Since the actual density curve of the residuals lies within the confidence envelope, we conclude that the normality assumption seems tenable.

## Studentized Residuals

Often researchers standardize the residuals before performing the assumption checking. This does not change any of the previous findings, in fact whether you use the raw residuals or the standardized residuals the scatterplot and density plot look identical, albeit with different scales. (Recall that standardizing is a linear function which does not change the shape of the distribution.)

$$z_{\epsilon} = \frac{\epsilon - 0}{\sigma_{\epsilon}}$$

(Remember the mean value of the residual at each  $X$  is 0, so that is why we subtract 0.) Unfortunately, we do not know what the value for  $\sigma_{\epsilon}$  (the standard error of the residuals) is, so we need to estimate it from the data. This adds uncertainty to the normal distribution and makes the distribution  $t$ -distributed.



$$t_{\epsilon} = \frac{\epsilon - 0}{\hat{\sigma}_{\epsilon}}$$

Since the  $t$ -distribution is also referred to as *Student's distribution*, this transformation of the residuals is called *studentizing*. What studentizing does for us is to put the residuals on a scale that uses the standard error. This allows us to judge whether particular residuals that look extreme (either highly positive or negative) are actually extreme or not. The studentized residuals are given in the augmented output as `.std.resid`.

```
# Plot the studentized residuals versus time spent on homework
ggplot(data = out_1, aes(x = homework, y = .std.resid)) +
  geom_point() +
  theme_bw() +
  geom_hline(yintercept = 0) +
  geom_hline(yintercept = c(-2, 2), linetype = "dashed") +
  xlab("Time spent on homework (in hours)") +
  ylab("Studentized residuals")

# Density plot of the studentized residuals
ggplot(data = out_1, aes(x = .std.resid)) +
  stat_watercolor_density(model = "normal") +
  stat_density(geom = "line") +
  theme_bw() +
  xlab("Studentized residuals") +
  ylab("Probability density")
```

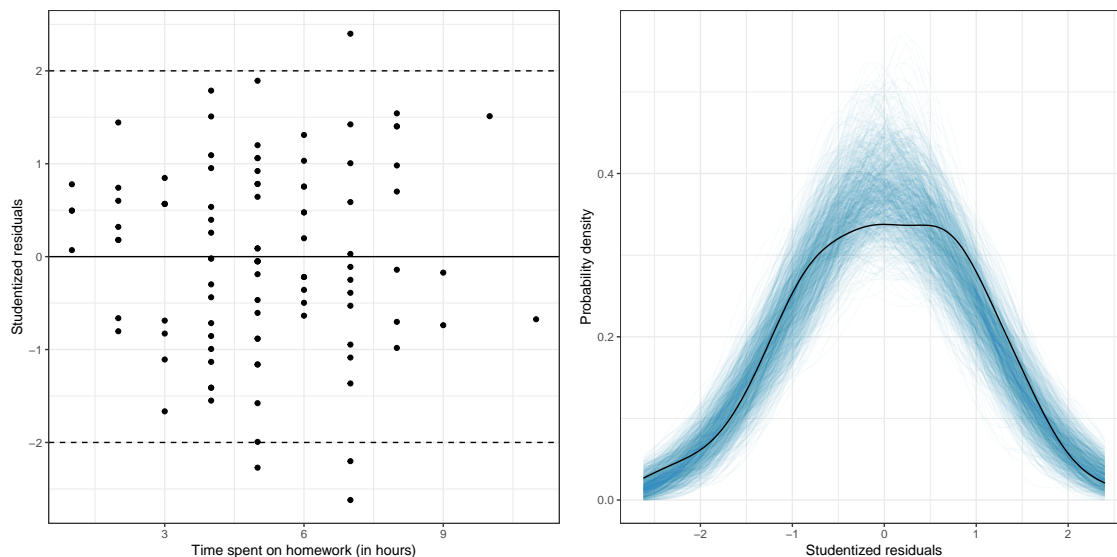


Figure 7. Scatterplot of the studentized residuals versus education level (left-hand side). Horizontal lines at two standard errors above and below  $Y = 0$  have been added to the plot to help identify observations that have a much higher or lower GPA than would be expected given their time spent on homework. A density plot of the studentized residuals is also shown (right-hand side). The bootstrapped confidence envelope (in blue) shows the sampling variation in density expected under the normality assumption.

The only thing that has changed between these plots and the previous plots of the raw residuals is the scale. However, now we can identify observations with extreme residuals, because we can make use of the fact that most of the residuals (~95%) should fall within two standard errors from the mean of 0. There are four students who

have residuals of more than two standard errors. Given that we have  $N = 100$  observations, it is not surprising to see four observations more extreme than two standard errors; remember we expect to see 5% just by random chance. If observations have really extreme residuals (e.g.,  $|t_e| > 3.5$ ), it is often worth a second look since these extreme observations are interesting and may point to something going on in the data.

We can also `filter()` the augmented data to find these observations and to determine the exact value of the studentized residuals. Recall that the vertical line (`|`) means “OR”.

```
out_1 %>%
  filter(.std.resid <= -2 | .std.resid >= 2)
```

# A tibble: 4 x 9

	gpa	homework	.fitted	.se.fit	.resid	.hat	.sigma	.cooksd	.std.resid
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	64	7	82.8	0.991	-18.8	0.0187	7.02	0.0655	-2.62
2	67	7	82.8	0.991	-15.8	0.0187	7.09	0.0462	-2.20
3	64	5	80.4	0.725	-16.4	0.0100	7.08	0.0261	-2.27
4	100	7	82.8	0.991	17.2	0.0187	7.06	0.0549	2.40

## Multiple Regression Assumptions

Recall that the model for a multiple regression (with two predictors) is a fitted plane that is composed of  $(x_1, x_2, y)$  ordered triples. Figure 8 visually shows the multiple regression model’s assumptions.

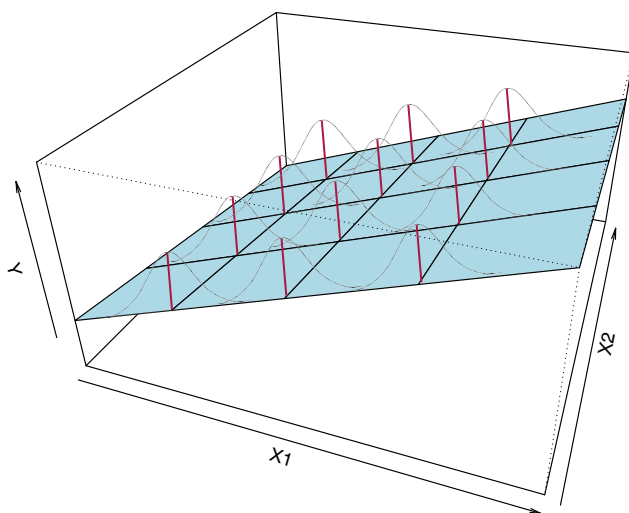


Figure 8. A visual depiction of the multiple regression model’s assumptions.

Now the  $Y$ -values, and thus the residuals, at each combination  $(x_1, x_2)$  are the conditional distributions we need to be thinking about.

The assumptions for the multiple regression model are similar to those for the simple model, namely,

- **Linearity:** The MEAN values of each combination  $(x_1, x_2)$  are linear in both the  $X_1$  and the  $X_2$  directions. The mean of each of the conditional distributions of the residuals is zero.
- **Independence:** Again, this is not shown in the figure. The assumption is that each residual value in a particular conditional distribution is independent from every other residual value in that same distribution.
- **Normality:** This assumption indicates that each of the conditional distributions of residuals is normally distributed.
- **Homoskedasticity:** The variance (or standard deviation) of all of the conditional distributions of residuals is exactly the same.

To evaluate these assumptions, we will create the exact same plots we created to evaluate the assumptions in the simple regression model, with one twist. Rather than creating the scatterplot by plotting the studentized residuals versus the predictor value, we will plot them against the FITTED values (i.e., the  $\hat{Y}_i$  values). The fitted values from a multiple regression represent the weighted combination of both predictors, and thus give us the appropriate conditioning when we examine the distributions. (Remember, we want to consider the distribution of residuals at each  $(x_1, x_2)$  combination.)

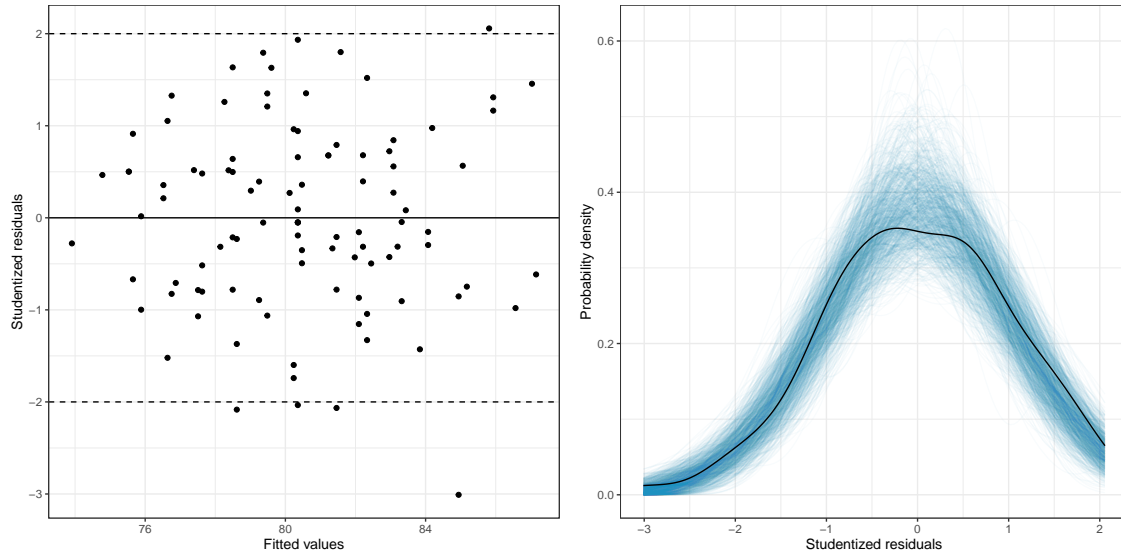
```
# Fit the multiple regression model
lm.2 = lm(gpa ~ 1 + homework + parent_ed, data = keith)

# Augment the model to obtain the fitted values and residuals
out_2 = augment(lm.2)
head(out_2)
```

```
# A tibble: 6 x 10
  gpa homework parent_ed .fitted .se.fit .resid .hat .sigma .cooks
  <dbl>   <dbl>   <dbl>   <dbl>   <dbl> <dbl> <dbl> <dbl>   <dbl>
1    78       2      13    76.5    1.29   1.48 0.0330  7.13 5.12e-4
2    79       6      14    81.3    0.783 -2.34 0.0122  7.12 4.54e-4
3    79       1      13    75.5    1.59   3.47 0.0500  7.12 4.41e-3
4    89       5      13    79.5    0.808   9.52 0.0130  7.06 8.01e-3
5    82       3      16    80.1    1.40   1.88 0.0390  7.13 9.88e-4
6    77       4      13    78.5    0.853 -1.50 0.0145  7.13 2.21e-4
# ... with 1 more variable: .std.resid <dbl>
```

```
# Plot the studentized residuals versus the fitted values
ggplot(data = out_2, aes(x = .fitted, y = .std.resid)) +
  geom_point() +
  theme_bw() +
  geom_hline(yintercept = 0) +
  geom_hline(yintercept = c(-2, 2), linetype = "dashed") +
  xlab("Fitted values") +
  ylab("Studentized residuals")

# Density plot of the studentized residuals
ggplot(data = out_2, aes(x = .std.resid)) +
  stat_watercolor_density(model = "normal") +
  stat_density(geom = "line") +
  theme_bw() +
  xlab("Studentized residuals") +
  ylab("Probability density")
```



*Figure 9.* Scatterplot of the studentized residuals versus the fitted values (left-hand side). Horizontal lines at two standard errors above and below  $Y = 0$  have been added to the plot to help identify observations that have a much higher or lower GPA than would be expected given their time spent on homework and parent level of education. A density plot of the studentized residuals is also shown (right-hand side). The bootstrapped confidence envelope (in blue) shows the sampling variation in density expected under the normality assumption.

The scatterplot shows random scatter around the  $Y = 0$  line which indicates that the linearity assumption seems tenable. The range of the studentized residuals at each fitted value also seem roughly the same indicating that the homoskedasticity assumption is also tenable. The density plot of the studentized residuals lies within the confidence envelope showing the random variation expected under the normal distribution, which suggests that the assumption of normality is tenable. Lastly, since the observations were randomly sampled we believe the independence assumption is satisfied.

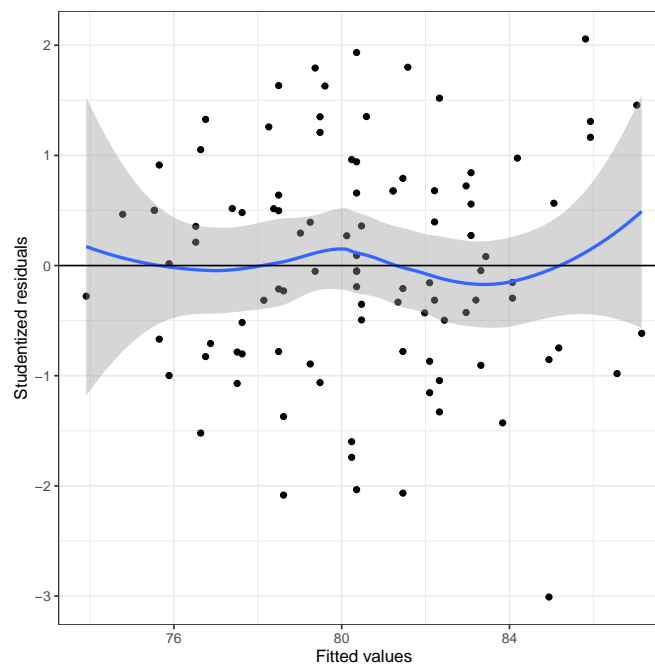
If any of the assumptions (aside from the independence assumption) do not seem reasonably satisfied, you can replot the residual plots based on the different simple regression models. (In this case we would look at the residuals versus time spent on homework and then the residuals versus parent education). This might help you identify if one, or both, of the predictors is the cause of the problem.

## Advanced Plotting: Loess Smooth to Help Evaluate Linearity

In the scatterplot of the studentized residuals (or raw residuals) versus the fitted values, we would expect that the average value of the residual at a given fitted value would be 0. The loess smoother helps us visualize the mean pattern in the actual data. We can then compare this to what would be expected (a constant mean pattern of 0) to evaluate the linearity assumption.

To add a loess smoother, we use the `geom_smooth()` function with the argument `method="loess"`. This will plot the loess line and also the confidence envelope around that loess line. This gives us an indication of the mean pattern in the data and its uncertainty. We would hope to see the line  $Y = 0$  (our expected pattern under linearity) encompassed in the uncertainty.

```
ggplot(data = out_2, aes(x = .fitted, y = .std.resid)) +  
  geom_point() +  
  geom_smooth(method = "loess") +  
  theme_bw() +  
  geom_hline(yintercept = 0) +  
  xlab("Fitted values") +  
  ylab("Studentized residuals")
```



*Figure 10.* Scatterplot of the studentized residuals versus the fitted values from a regression model using time spent on homework and parent education level to predict GPA. A horizontal line at  $Y = 0$  shows the expected mean residual under the linearity assumption. The loess line (blue) and uncertainty bands (grey shaded area) are also displayed.

Note that the loess line shows some deviation from the  $Y = 0$  line, however this is likely just due to sampling variation as the line  $Y = 0$  is encompassed in the confidence envelope. Because of this we would suggest that the linearity assumption is tenable.

## Advanced Plotting: Identify Observations with Extreme Residuals

It can be useful to identify particular observations in the residual plots directly. This can be useful as you explore the plots, and also to create plots for publications in which you wish to highlight particular cases. In the plot below, we will use the fortified data from the initial simple regression model we fitted to identify the cases in the scatterplot.

Rather than plotting points (`geom_point()`) for each observation, we can plot text for each observation using `geom_text()`. For example, you might imagine writing the name of each student in place of their point on the scatterplot. To do this, we first need to create an ID variable in the augmented data, then use `geom_text()` rather than `geom_point()` in the ggplot syntax. In the `geom_text()` function we will set `label=` to the newly created ID variable, and since it is a variable in the data set, we will put that in an `aes()` function.

Since the original data set does not include an ID variable (e.g., names), we will use the row number from the original data as the ID. In other words the student in the first row will have an ID of 1, etc.

```
# Create ID variable in the augmented data
out_2 = out_2 %>%
  mutate(id = row.names(keith))

head(out_2)
```

```
# A tibble: 6 x 11
   gpa homework parent_ed .fitted .se.fit .resid .hat .sigma .cooksd
<dbl>   <dbl>    <dbl>   <dbl>   <dbl> <dbl> <dbl> <dbl>   <dbl>
1    78      2      13    76.5    1.29  1.48 0.0330  7.13 5.12e-4
2    79      6      14    81.3    0.783 -2.34 0.0122  7.12 4.54e-4
3    79      1      13    75.5    1.59  3.47 0.0500  7.12 4.41e-3
4    89      5      13    79.5    0.808  9.52 0.0130  7.06 8.01e-3
5    82      3      16    80.1    1.40  1.88 0.0390  7.13 9.88e-4
6    77      4      13    78.5    0.853 -1.50 0.0145  7.13 2.21e-4
# ... with 2 more variables: .std.resid <dbl>, id <chr>
```

```
# Plot the id variable as text rather than points in the scatterplot
ggplot(data = out_2, aes(x = .fitted, y = .std.resid)) +
  geom_text(aes(label = id), size = 4) +
  theme_bw() +
  geom_hline(yintercept = 0) +
  geom_hline(yintercept = -2, linetype = "dotted") +
  geom_hline(yintercept = 2, linetype = "dotted") +
  xlab("Fitted values") +
  ylab("Studentized residuals")
```

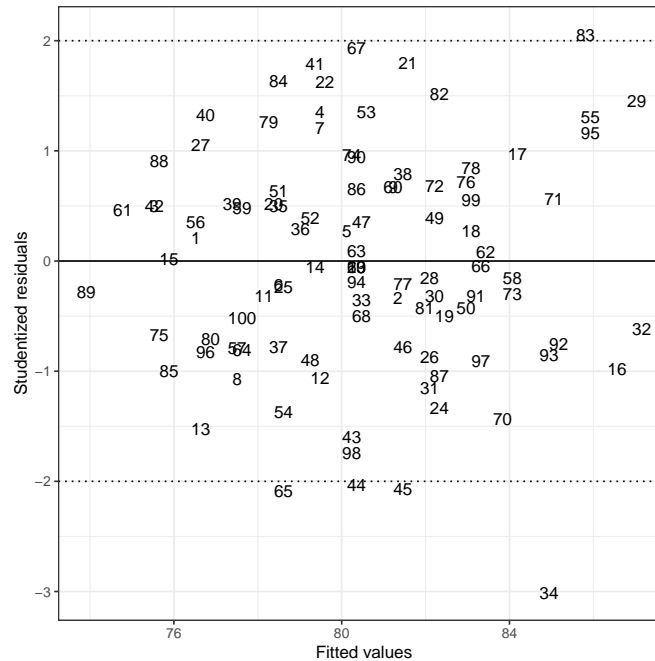


Figure 11. Scatterplot of the studentized residuals versus the fitted values from a regression model using time spent on homework and parent education level to predict GPA. A horizontal line at  $Y = 0$  shows the expected mean residual under the linearity assumption. The values plotted indicate the students' row numbers in the data.

We can also plot points for some students and ID label for other students. For example, suppose we wanted to give the ID number for only those students with a studentized residual that was less than  $-2$  or greater than  $2$ , and plot a point otherwise. To do this, we would create the ID variable in the augmented data (which we have already done), then split the dataset into two datasets: one for those employees with extreme residuals and one for those that have a non-extreme residual. Then we will call `geom_point()` for those in the non-extreme data set, and `geom_text()` for those in the extreme set. We do this by including a `data=` argument in one of those functions to call a different dataset.

```
# Create different data sets for the extreme and non-extreme observations
extreme = out_2 %>%
  filter(.std.resid <= -2 | .std.resid >= 2)

nonextreme = out_2 %>%
  filter(.std.resid > -2 & .std.resid < 2)

# Plot using text for the extreme observations and points for the non-extreme
ggplot(data = extreme, aes(x = .fitted, y = .std.resid)) +
  geom_text(aes(label = id), size = 4, color = "red") +
  geom_point(data = nonextreme) +
  theme_bw() +
  geom_hline(yintercept = 0) +
  geom_hline(yintercept = -2, linetype = "dotted") +
  geom_hline(yintercept = 2, linetype = "dotted") +
  xlab("Fitted values") +
  ylab("Studentized residuals")
```

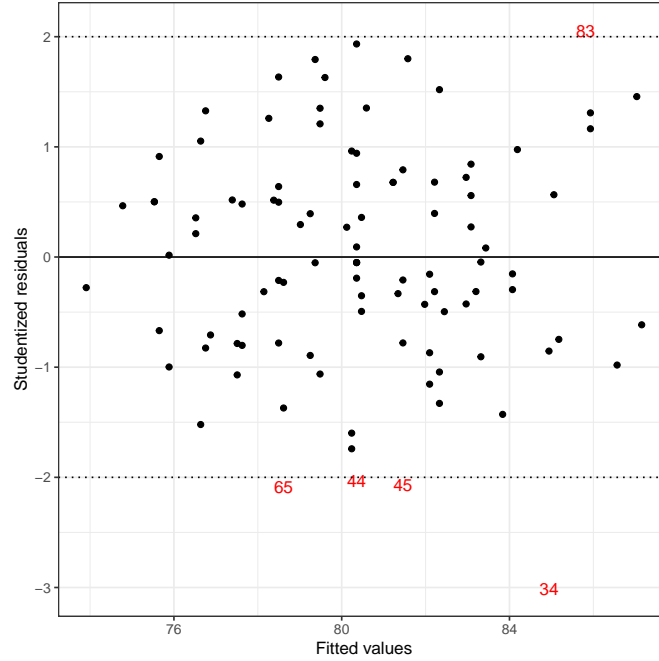


Figure 12. Scatterplot of the studentized residuals versus the fitted values from a regression model using time spent on homework and parent education level to predict GPA. A horizontal line at  $Y = 0$  shows the expected mean residual under the linearity assumption. Students with studentized residuals more than two standard errors from 0 are also identified by their row number.

## Regression Model

To date we have been writing the regression model as a mathematical expression of the relationship between some outcome ( $Y$ ) and a set of predictors ( $X_1, X_2, \dots, X_k$ ), namely as,

$$Y_i = \beta_0 + \beta_1(X_{1i}) + \beta_2(X_{2i}) + \dots + \beta_k(X_{ki}) + \epsilon_i$$

This is partially correct. A statistical model needs to represent the data generating process, which also embodies the set of underlying assumptions. This implies that the regression model encompasses both the mathematical relationship and the underlying assumptions:

$$Y_i = \beta_0 + \beta_1(X_{1i}) + \beta_2(X_{2i}) + \dots + \beta_k(X_{ki}) + \epsilon_i$$

where  $\epsilon_{i|X} \stackrel{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma^2)$