

## Lecture 2.2 - Simple regression activity

Student

2025-11-06



### Simple regression

#### Planning

#### Expectations

In this exercise, we are going to try and predict calories

1. Pick a predictor variable that you think will affect calories

## Saturated fat

Write down your expected relationship between the two variables – specify what you think the correlation will be (high, medium, low, positive, negative)

I expect there to be a positive, somewhat linear, and moderate relationship between the variables

2. Create a scatterplot of the relationship between the two variables

You can create a scatterplot (replacing the text `<...>` with your data names) with:

```
ggplot(<dataset>, aes(x=<variable1>, y=<variable2>)) + geom_point()
```

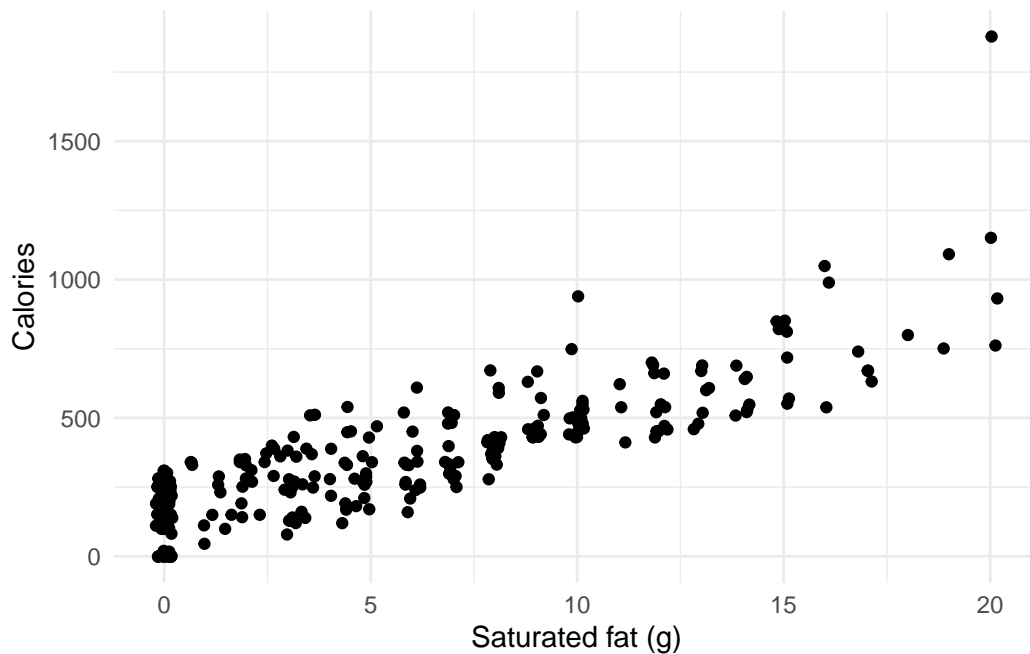


Figure 1: Relationship between calories and saturated fat

- Make a high-quality scatterplot of your two variables.
- Write down a brief description of the relationship between the two, including your prediction of the correlation of the two and what values you expect the slope and intercept to be of the best fit line

There seems to be a relatively strong relationship between the two. I would guess that the intercept is around 200 and the slope is about 50

## Investigation

3. Calculate the correlation between the two variables and record it.

- Was it stronger or weaker than you expected?

Remember, the code for correlation is:

```
cor(<dataset>$<variable1>, <dataset>$<variable2>, use="complete.obs")
```

Replace <dataset> with your dataset's name and <variable> with the variable's name that you are analyzing.

Table 1: Correlation between Calories and Total fat

Variable1	Variable2	Correlation
Calories	Saturated fat (g)	0.846

As expected, the relationship is strong.

4. Now let's create some lines through the data. Below is sample code you will need to copy and paste into a code block.

```
scatterplot <- ggplot(<dataset>, aes(x=<predictor.variable>, y=<response.variable>))  
  + geom_point()
```

```
scatterplot  
  + geom_abline(slope=<your.slope.estimate>, intercept=<your.intercept.estimate>)
```

Remember to replace the parts in <> with your own information.

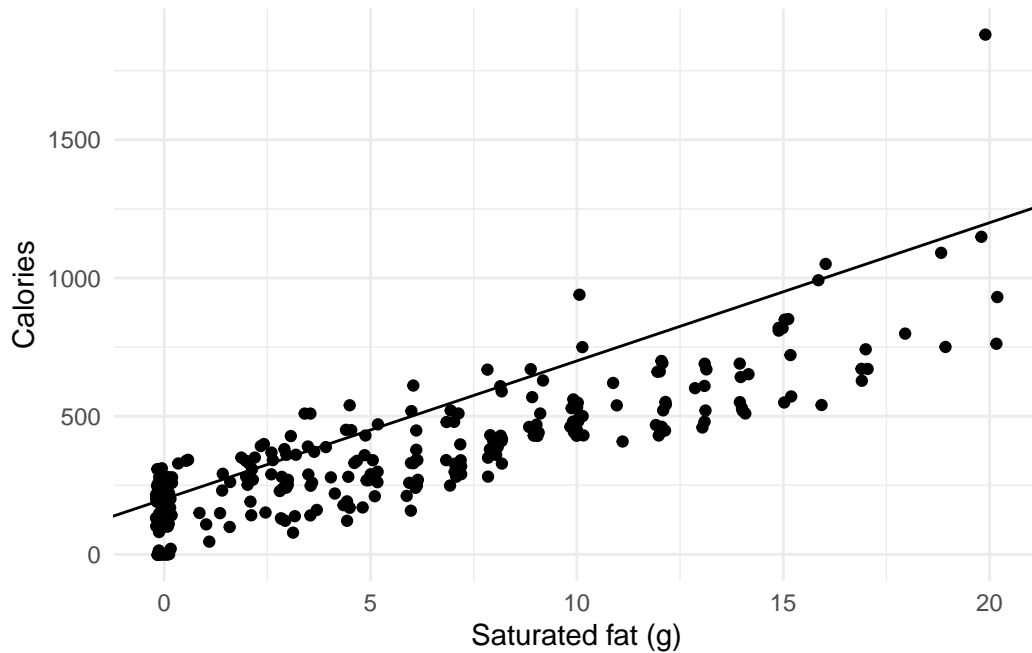


Figure 2: Estimate of best fit line for Calories and Saturated fat

- How well do you feel your best fit line fits the data?

Slope line is probably a little high but looks to be in the right ballpark

5. Next run a regression to estimate the linear relationship between the two variables. We can do this with:

```
lmmodel <- lm(<response.variable> ~ <predictor.variable>, data=<dataset>)
summary(lmmodel)
```

Remember to replace `<response.variable>`, `<predictor.variable>` and `<dataset>` with the data and variables you are using.

Call:

```
lm(formula = Calories ~ Saturated.Fat, data = mcdonalds)
```

Residuals:

Min	1Q	Median	3Q	Max
-209.73	-106.61	-21.59	81.31	977.57

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	138.924	12.033	11.54	<2e-16 ***
Saturated.Fat	38.175	1.501	25.44	<2e-16 ***

---

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

Residual standard error: 128.5 on 258 degrees of freedom

Multiple R-squared: 0.715, Adjusted R-squared: 0.7139

F-statistic: 647.2 on 1 and 258 DF, p-value: < 2.2e-16

- How close to your line was the regression line? If it was different, why do you think that happened?

Not too far off, mine was just a guess

- Interpret your results – for every change in your predictor variable, how much does your response variable change?

For every gram increase in saturated fat, the predicted calories increase by 38.

- Compared to your expectations, is the predicted change a lot or a little?

It's within roughly the range expected

## Checking your work

7. For the final step of part I, add a regression line to your plot. Do this by inserting the following commands into a new code block:

```
scatterplot <- ggplot(<dataset>, aes(x=<predictor.variable>, y=<response.variable>))  
  + geom_point()
```

```
scatterplot  
  + geom_abline(slope=<regression.slope.estimate>,  
                intercept=<regression.intercept.estimate>)
```

Instead of using your estimate from before, this time enter the slope and intercept generated by the regression in the previous step.

```
ggplot(mcdonalds, aes(x=Saturated.Fat, y=Calories)) +  
  geom_jitter() +  
  labs(x="Saturated fat (g)", y="Calories") +  
  geom_abline(slope=38.175, intercept=138.924)
```

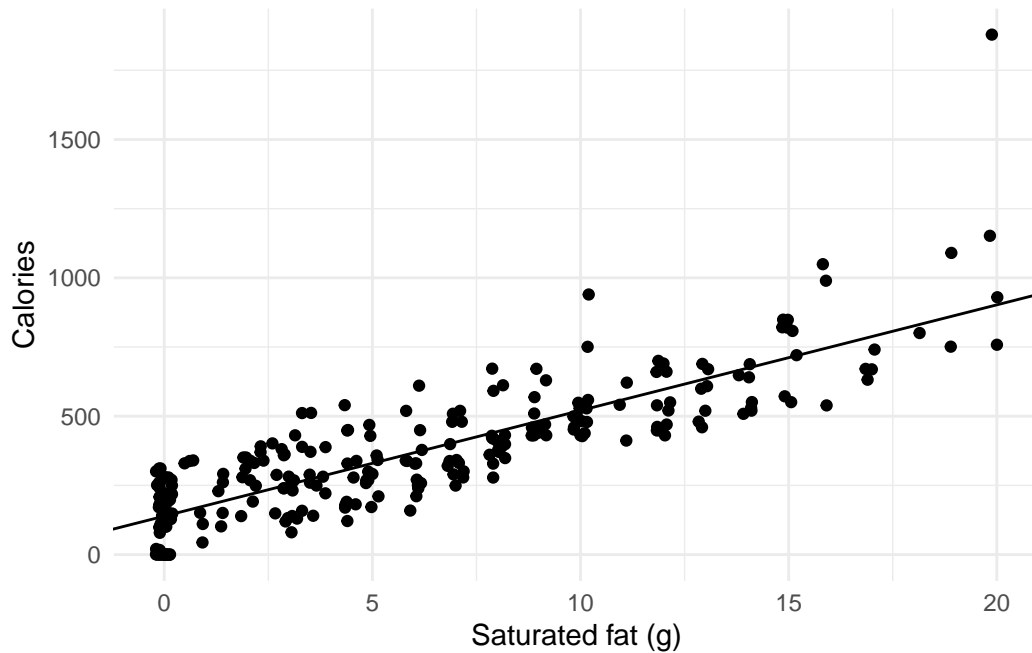


Figure 3: OLS line of best fit line for Calories and Saturated fat

- How close was your predicted best fit line compared to the regression line R generated?  
If they were different, why were they different?

Fairly close

- Would either of your variables benefit from a transformation?

In this case, it seems no

## Simple regression extension

### Calculating residuals

8. Make a small table using the Quarto (manually make a table in visual mode) built-in table function. Make a table with five columns and four rows. Then pick any three observations in the dataset and, in the table, record the following information:

- Item name
- Item's calories
- Item's observed value for calories
- The predicted value of calories based on the regression model

- The size of the residual.

Table 2: Residual calculations for Calories and Total fat model

Item name	Actual saturated fat	Actual calories	Predicted calories	Residual size
Chocolate Chip Cookie	3.5	160	272.537	-112.537
Coca-Cola Classic (Small)	0	140	138.924	1.07
Iced Mocha (Medium)	8	350	444.324	94.324

### Residual calculations for Calories and Total fat model

How large were your residuals? Did the size of the residuals indicate to you that the regression line is a good fit or not?

Fairly large, though I think we got unlucky on the random item selection

### Residual plotting

#### 9. Make a residual plot & standard deviation

Fortunately, RStudio can easily make a residual plot. When you run a regression, within the stored regression RStudio stores the values of  $\beta$ .

Note that you will need to load **broom** library for this code to work. You will be using the `augment()` function that adds the residuals to a dataset for easy display and manipulation.

Note that in the lecture we viewed the predictor variable on the  $x$  axis of the residual plot. Here we are switching to viewing the response variable on the  $x$  axis here. If you think about it, the two plots are equivalent, simply rotated. But once we switch to thinking in higher dimensions, the only plot that makes sense is the residual plot with the response variable plotted.

```
# Save the model
<modelname> <- lm(<predictor.variable> ~ <response.variable>, data=<dataset>)
# Create the residuals database
<modelname>.augmented <- augment(<modelname>, <dataset>)

# Residual histogram plot
```

```
ggplot(<modelname>.augmented, aes(x=.resid)) +
  geom_histogram(fill="blue4") +
  labs(x="Residuals", y="Count")

# Residual scatterplot
ggplot(<modelname>.augmented, aes(<response variable>, .resid)) +
  geom_point() +
  geom_hline(yintercept = 0, color = "blue", linetype='dashed') +
  labs(y = "Residuals", x="<name of response variable>")
```

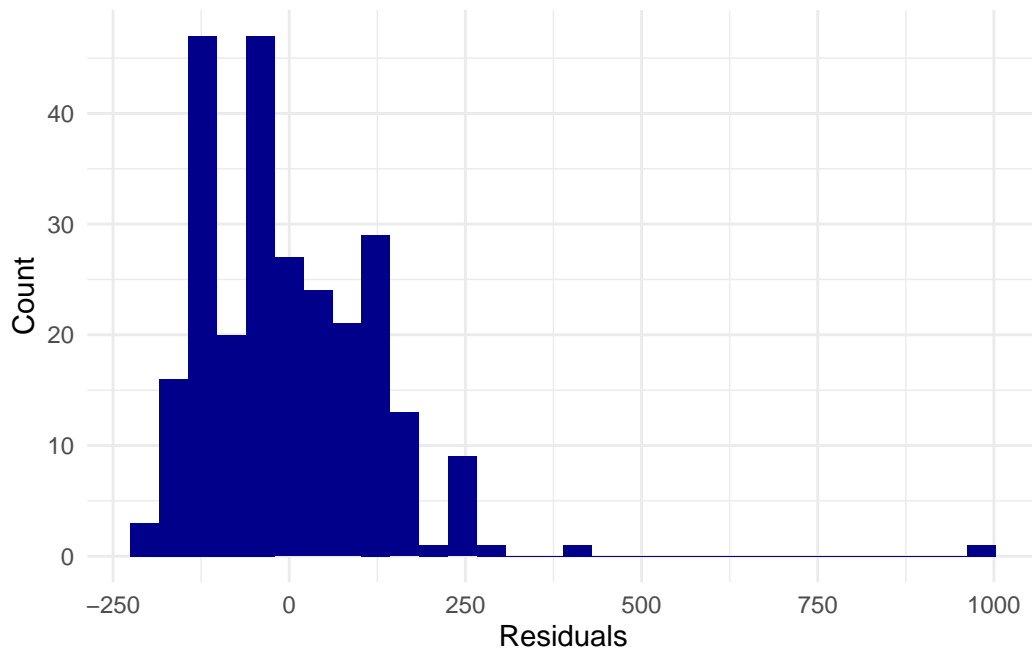


Figure 4: Histogram of the residuals

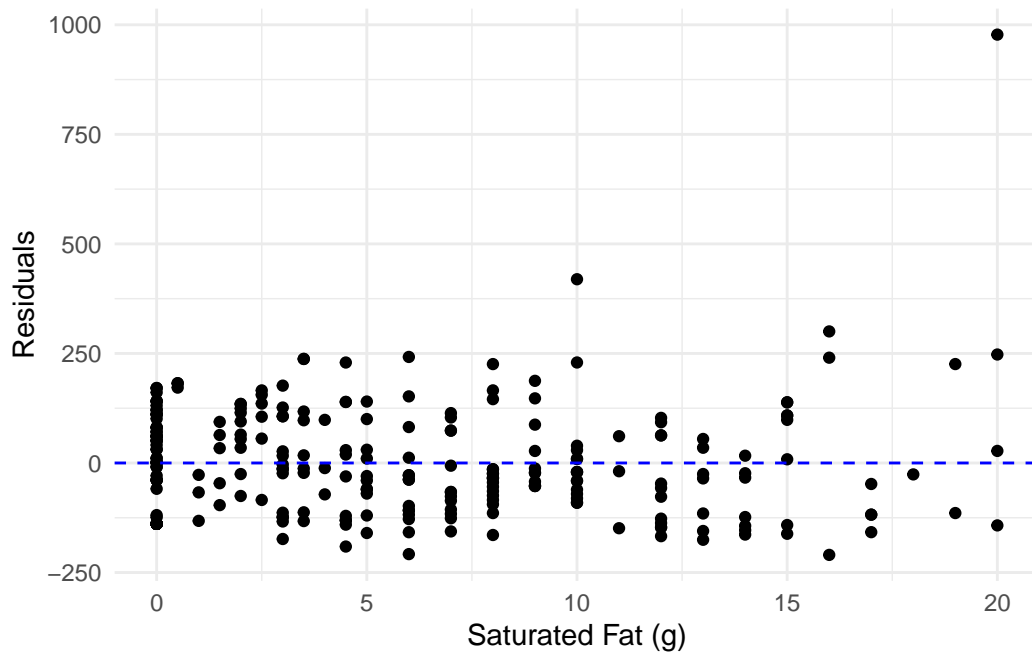


Figure 5: Scatterplot of the residuals

Residual information for Calories and Saturated Fat model

- What can you conclude from the residual plots and the standard deviation calculation of the residuals? Does it indicate a good model fit or not?

Overall, the model seems to be very good fit with one outlier

## Model fit

10. Interpret the R squared

What can you conclude about your model based on the R squared?

```
summary(calmodel)$r.squared
```

```
[1] 0.7149779
```

The R squared is quite high, which fits with the scatterplot created earlier - the model is a relatively good fit for the data.

11. Which observations are obvious outliers? How does your line of best fit change if you exclude some of the outliers?

```
mcdonalds.nooutliers <- mcdonalds %>%  
  filter(Calories < 1400)  
  
summary(lm(Calories ~ Saturated.Fat, data=mcdonalds.nooutliers))
```

Call:

```
lm(formula = Calories ~ Saturated.Fat, data = mcdonalds.nooutliers)
```

Residuals:

Min	1Q	Median	3Q	Max
-204.11	-93.50	-13.48	79.64	430.88

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	146.602	10.607	13.82	<2e-16 ***
Saturated.Fat	36.252	1.336	27.13	<2e-16 ***

---

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

Residual standard error: 112.9 on 257 degrees of freedom

Multiple R-squared: 0.7412, Adjusted R-squared: 0.7402

F-statistic: 736 on 1 and 257 DF, p-value: < 2.2e-16

The coefficients hardly change when we remove the 40 piece chicken nuggets - it was a high leverage but low influence point.

If you are finished with this analysis, conduct the exercise again for a different predictor variable.