

Lesson 6 - Linear Regression, Multiple Linear Regression, ANOVA, ANCOVA: Choosing the Best Model for the Job

Contents

Answering questions with data	2
Our Dataset	3
T-tests	6
How we Evaluate which Model to Use	8
Assumptions of general linear models	9
The Linear Models	10
Simple linear regression	12
Multiple linear regression	14
Adding powers of a variable (polynomial regression)	14
Adding extra variables to our model	15
Interaction terms	18
One-way analysis of variance (ANOVA)	19
Dummy Variables	20
Multiple test correction	22
Multi-way analysis of variance (ANOVA)	24
Analysis of covariance (ANCOVA)	26
Review: Models we used today	30
Appendix	32
Prediction	32
Assessing the performance of the model (feedback)	34
Checking Residuals	34
QQ-plots	37
Next Steps (or When Assumptions Fail)	41

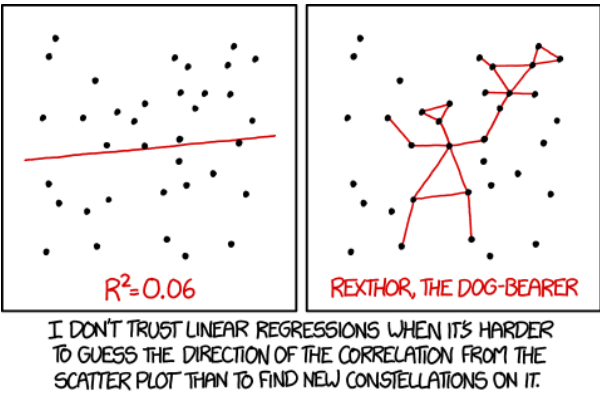


Figure 1: xkcd

Objective: At the end of this session you will be able to perform simple and multiple linear regression, one- and multiway analysis of variance (ANOVA) and analysis of covariance (ANCOVA). You will be able to interpret the statistics that come out of this model, be cognizant of the assumptions the model makes, and use an F-test to select the best model for the job.

Highlighting

grey background - a package, function, code or command

italics - an important term or concept

bold - heading or a term that is being defined

blue text - named or unnamed hyperlink

Packages Used in This Lesson

The following packages are used in this lesson:

`tidyverse` (`ggplot2`, `tidyr`, `dplyr`)

`limma`

`gee`

`multcomp`

`broom`

Please install and load these packages for the lesson. In this document I will load each package separately, but I will not be reminding you to install the package. Remember: these packages may be from CRAN OR Bioconductor.

Data Files Used in This Lesson

-SISG-Data-cholesterol.txt

These files can be downloaded at https://github.com/eacton/CAGEF/tree/master/Lesson_6/data. Right-click on the filename and select ‘Save Link As...’ to save the file locally. The files should be saved in the same folder you plan on using for your R script for this lesson.

Or click on the blue hyperlink at the start of the README.md at https://github.com/eacton/CAGEF/tree/master/Lesson_6 to download the entire folder at DownGit.

Answering questions with data

In order to work with our data we need to be able to answer some basic questions.

- How do we describe our data?
- How do we test our hypotheses (what model do we use)?
- How do we test our assumptions about the model are using?
- How do we compare models?
- How do we make a prediction with new values?

These are all really important questions that we may or may not think about as we try to dive in and get our answer as quickly as possible. Today we are going to slow down a bit and think about our data and our models.

Load the packages!

```
library(tidyverse)
library(limma)
library(gee)
library(multcomp)
library(broom)
```

Our Dataset

The dataset we will use for this lesson is from the Summer Institute in Statistical Genetics (SISG) at the University of Washington's course in Regression and Analysis of Variance by Lurdes Inoue. This lesson uses a lot of material from the SISG 2016 course as well as conceptual material from Ben Bolker. I like this dataset because it has a number of categorical and continuous variables, which allows us to use the same dataset for many models. Also, the variables are familiar (age, BMI, gender, cholesterol), which makes data interpretation easier while we are in the learning stage.

Read the data in and take a look at the structure.

```
cholesterol <- read.delim("data/SISG-Data-cholesterol.txt", sep = " ", header = TRUE)

str(cholesterol)
```

```
## 'data.frame':   400 obs. of  9 variables:
## $ ID           : int  1 2 3 4 5 6 7 8 9 10 ...
## $ sex          : int  1 1 0 0 1 1 0 0 0 0 ...
## $ age          : int  74 51 64 34 52 39 79 38 52 58 ...
## $ chol         : int  215 204 205 182 175 176 159 169 175 189 ...
## $ BMI          : num  26.2 24.7 24.2 23.8 34.1 22.7 22.9 24.9 20.4 22 ...
## $ TG           : int  367 150 213 111 328 53 274 137 125 209 ...
## $ rs174548     : int  1 2 0 1 0 0 2 1 0 0 ...
## $ rs4775401    : int  2 1 1 1 0 2 1 1 1 1 ...
## $ APOE         : int  4 4 4 1 1 4 1 1 4 5 ...
```

This dataset is looking at genetic variants (single nucleotide polymorphisms (SNPs)) and their relationship to cholesterol (chol) and triglycerides (TG) for 3 genes: rs174548 (FADS1 - an enzyme in fatty acid unsaturation), rs4775401 (a candidate SNP), and APOE (a major apolipoprotein important for Alzheimer's disease).

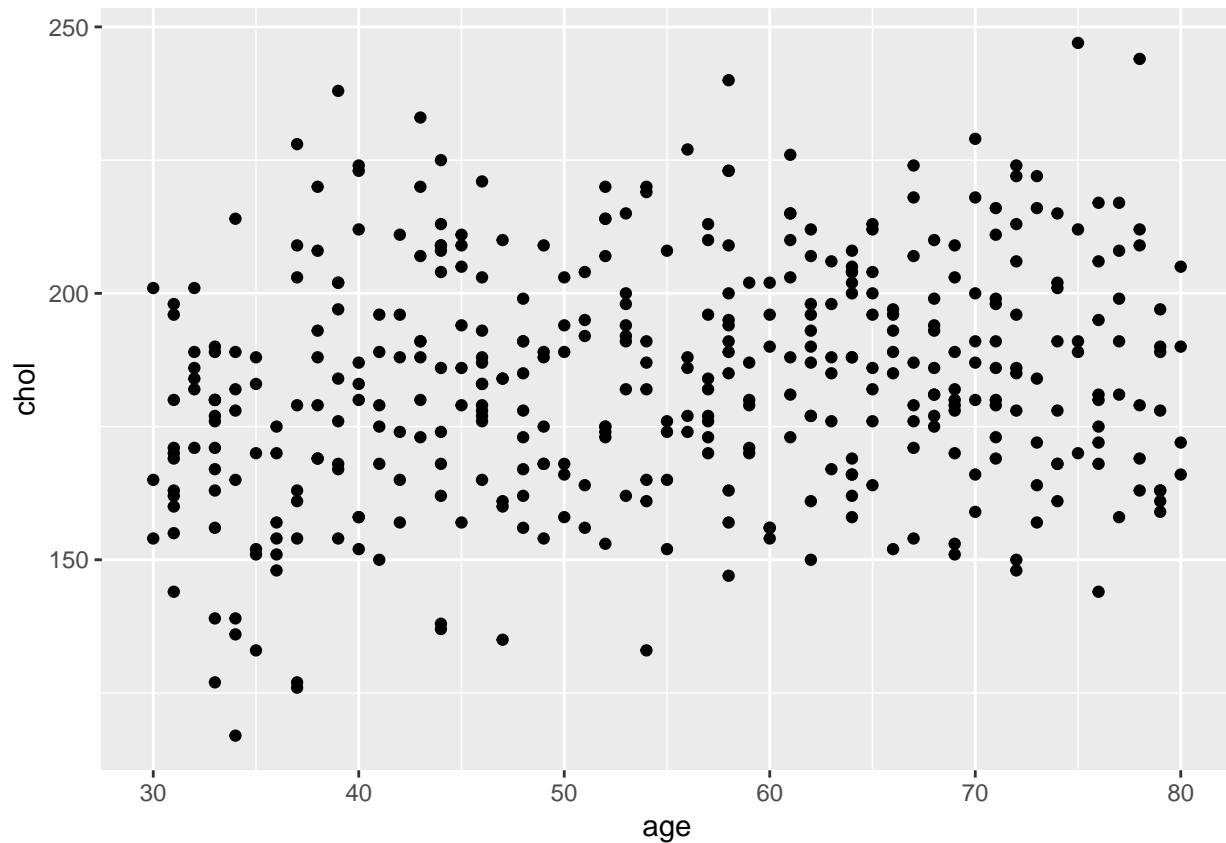
Note that categorical variables have been encoded. Sex is 0 and 1 instead of 'Male' and 'Female'. rs174548 has 3 possible nucleotide combinations, "C/C", "C/G", and "G/G" which have been encoded as 0, 1, and 2. Similarly rs4775401 has 0, 1, and 2 representing SNPs. APOE has 6 variants labelled starting at 1 (1-6).

We are ultimately interested in the relationship between the above genetic variants and cholesterol, while controlling for factors such as age and sex. But let's get our feet wet by starting with the easier **question: is there an association between mean serum cholesterol and age?**

For this question cholesterol is the **dependent variable**, or the variable being measured. Age is the **independent variable** that we are changing to determine the effect on cholesterol.

It is always, always, always, a good idea to make an 'exploratory' plot of your data and get an idea of what its distribution looks like. We can start with a simple scatterplot of age and cholesterol.

```
ggplot(cholesterol, aes(age, chol)) +
  geom_point()
```



We could also describe our data with some basic summary statistics such as the mean, median, mode, min, max, standard deviation, and variance. R does not have a mode function relating to statistics, but we can figure it out for ourselves.

```
cholesterol %>%
  summarize(mean = mean(chol), median = median(chol), min = min(chol),
            max = max(chol), sd = sd(chol), variance = var(chol))
```

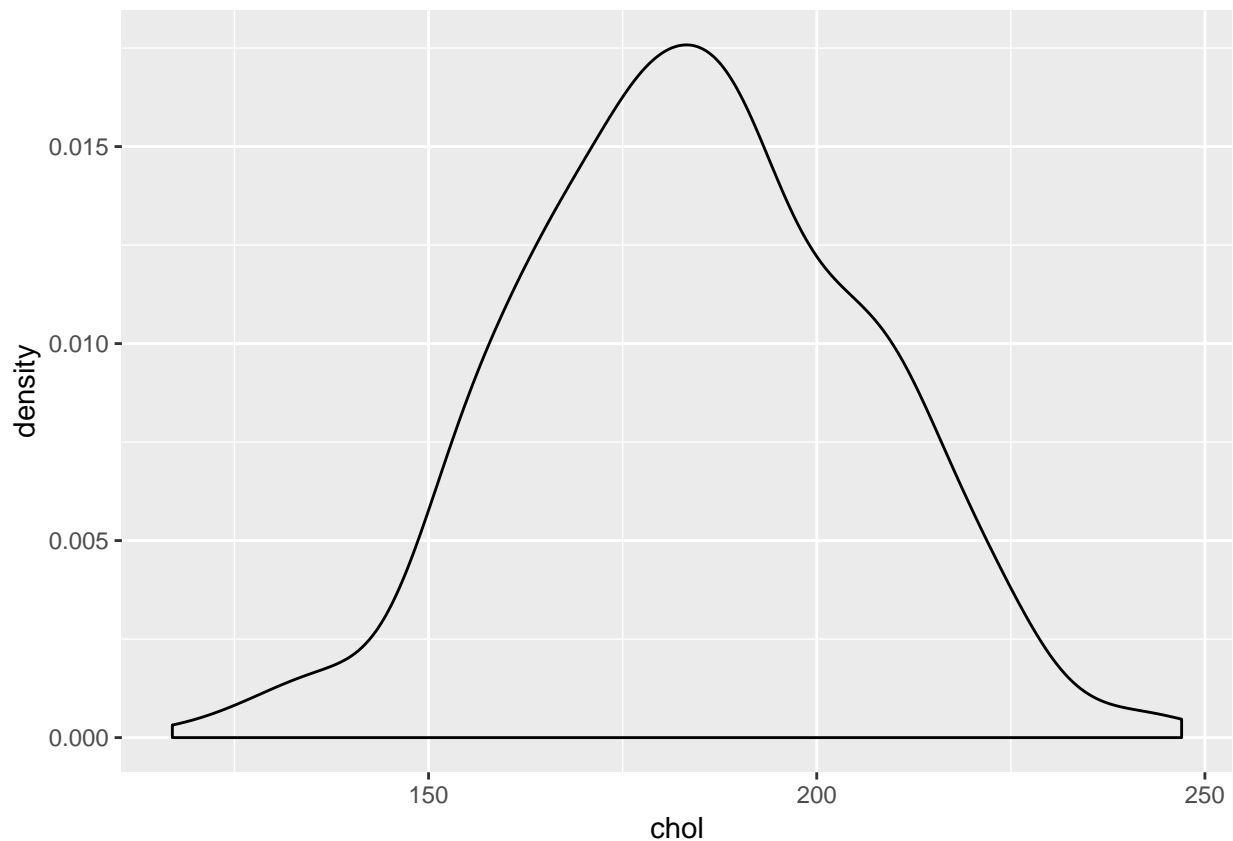
```
##      mean median min max      sd variance
## 1 183.915    184 117 247 22.11777 489.1958
```

```
mode <- sort(table(cholesterol$chol), decreasing = TRUE)
mode[1]
```

```
## 191
## 12
```

Our mean, median and mode are not that different, and so our data is not skewed in either direction. We can also prove this to ourselves by making a quick density plot.

```
ggplot(cholesterol, aes(chol)) +
  geom_density()
```



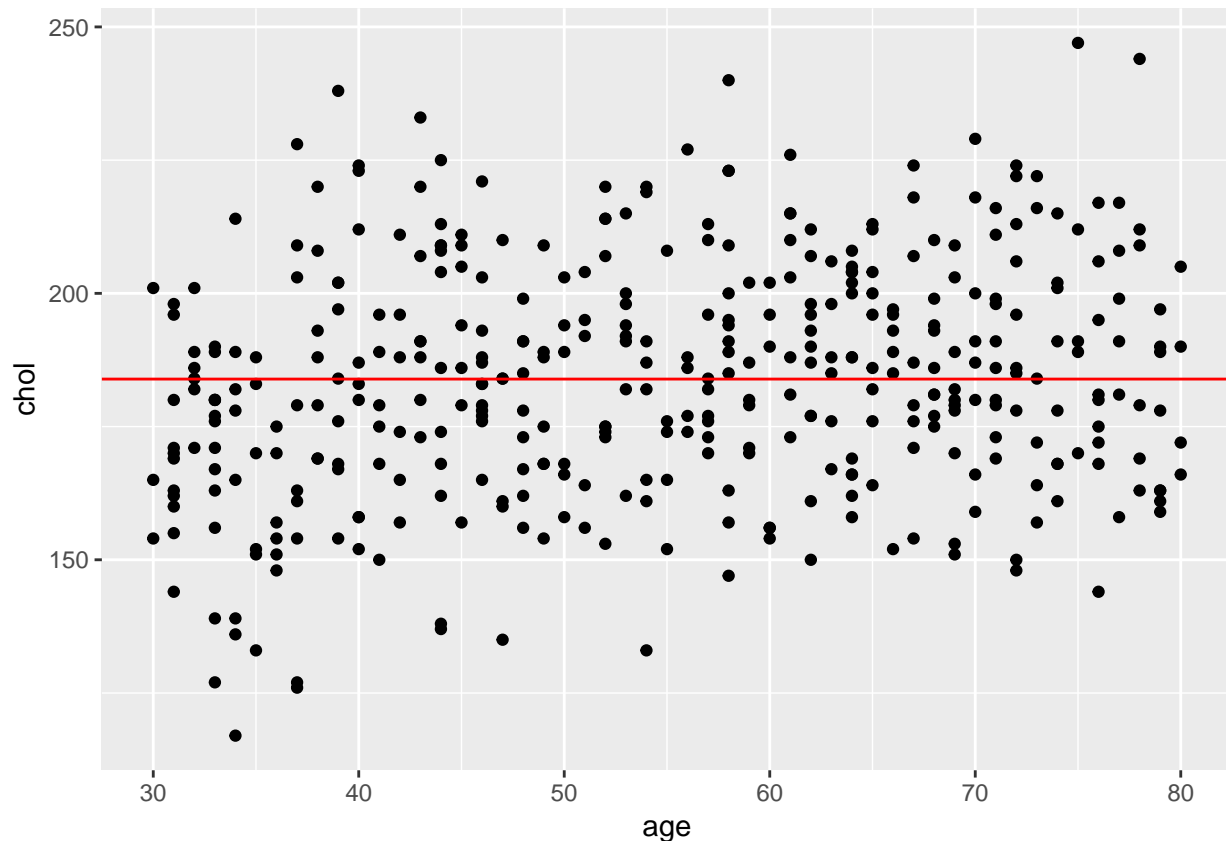
The `quantile()` function will also give us a good sense of the range and distribution of our data.

```
quantile(cholesterol$chol)
```

```
##      0%      25%      50%      75%     100%
## 117.00 168.00 184.00 199.25 247.00
```

Going back to our question of if age is related to cholesterol, let's add the mean cholesterol to our plot for reference. This is done by adding `geom_hline` and specifying the value for the 'yintercept'.

```
ggplot(cholesterol, aes(age, chol)) +
  geom_point() +
  geom_hline(yintercept = mean(cholesterol$chol), color = "red")
```



It looks like the mean might increase with age, but how do we test this?

T-tests

T-tests are a simple statistical tool let us to compare the means between groups. We don't currently have age groups, but we can make them. One way to do this is to use our `dplyr` skills to create a new column 'age_group'. The data can be split at 55 years-old (the midpoint of age in our data).

We can use an if/else statement (the `ifelse` function) to test: is age greater than 55? If the answer is 'yes' the value is 1 and if the answer is 'no' the value is 0. We can take a quick look at our dataset to make sure this worked.

```
cholesterol <- cholesterol %>% mutate(age_group = ifelse(test = cholesterol$age > 55, yes = 1, no = 0))
```

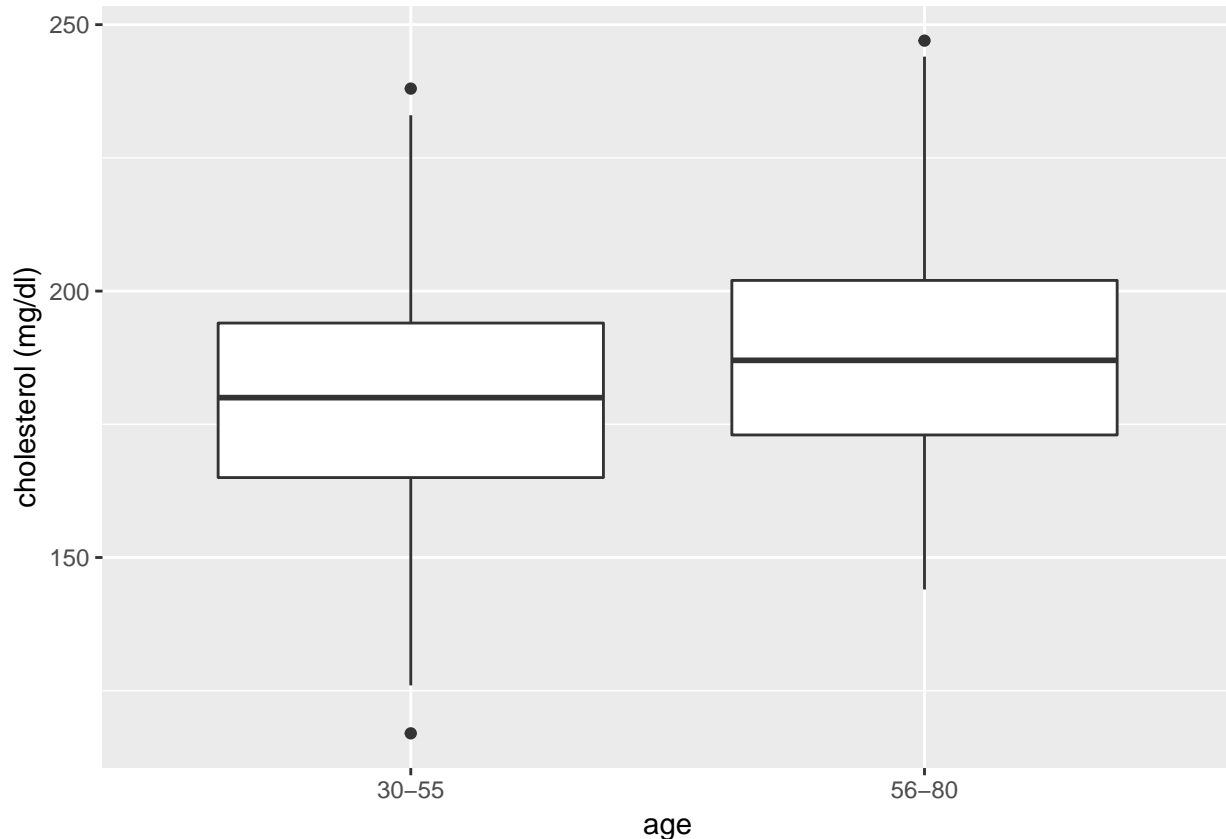
```
str(cholesterol)
```

```
## 'data.frame': 400 obs. of 10 variables:
## $ ID : int 1 2 3 4 5 6 7 8 9 10 ...
## $ sex : int 1 1 0 0 1 1 0 0 0 0 ...
## $ age : int 74 51 64 34 52 39 79 38 52 58 ...
## $ chol : int 215 204 205 182 175 176 159 169 175 189 ...
## $ BMI : num 26.2 24.7 24.2 23.8 34.1 22.7 22.9 24.9 20.4 22 ...
## $ TG : int 367 150 213 111 328 53 274 137 125 209 ...
## $ rs174548 : int 1 2 0 1 0 0 2 1 0 0 ...
## $ rs4775401 : int 2 1 1 1 0 2 1 1 1 1 ...
## $ APOE : int 4 4 4 1 1 4 1 1 4 5 ...
## $ age_group: num 1 0 1 0 0 0 1 0 0 1 ...
```

We can now use a boxplot to look at the distribution of cholesterol for our 2 groups.

Boxplots are a great way to visualize summary statistics for your data. As a reminder, the thick line in

```
ggplot(cholesterol, aes(factor(age_group),chol)) +  
  geom_boxplot() +  
  scale_x_discrete(labels = c("30-55", "56-80")) +  
  xlab("age") +  
  ylab("cholesterol (mg/dl)")
```



There seems to be a lot of overlap in our cholesterol values. How do we tell if the means are truly different?

Let's think about this a little more explicitly:

The **null hypothesis** is that there is no difference in the sample means between our groups.

An **alternative hypothesis** is that there is a difference between the means (2-sided test), or that the difference in means is greater or lesser than zero (1-sided test).

α is our p-value, and μ is the population mean, k is our sample mean. Remember that we are *estimating* the true population mean using the sample that have. Our **p-value** is the probability of finding our observed value by chance given that the null hypothesis is true.

We will use a simple student's t-test to test the alternative hypothesis that the true difference in means is not equal to 0.

The `t.test` function takes as input the variables on which we are performing the test (in vector format), the type of t-test being performed, and the confidence interval. The **confidence interval** is the interval that will cover the true parameter $x\%$ of the time. In the image above the confidence interval covers the pink area, $(1-\alpha)$.

You can alternatively enter your variables in a **formula**, in this case $y \sim x$. The \sim in r language is used

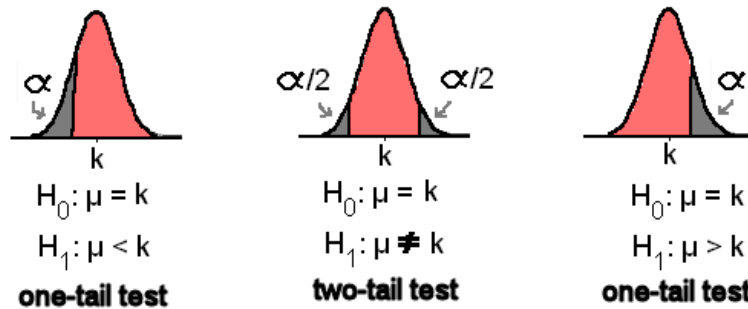


Figure 2:

to separate the left and right sides of a formula. (You can run a 1-sided t-test by specifying `alternative = 'greater'` or `alternative = 'less'`). In this case, `alternative = 'two-sided'` and `conf.level = 0.95` are the default parameters and only included for clarity. For now we are assuming that equal variance is true.

```
t.test(x= cholesterol$age_group, y = cholesterol$chol, alternative = "two.sided", conf.level = 0.95, var.equal = TRUE)

#is equivalent to
t.test(formula = cholesterol$chol ~ cholesterol$age_group, var.equal = TRUE)
```

```
##
## Two Sample t-test
##
## data:  cholesterol$chol by cholesterol$age_group
## t = -3.6349, df = 398, p-value = 0.0003146
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -12.202574 -3.636122
## sample estimates:
## mean in group 0 mean in group 1
##      179.9751      187.8945
```

Interpretation

Our output tells us the mean cholesterol for those aged 30-55 is 180 mg/dl and the mean for those aged 56-80 is 188 mg/dl. The difference in means is significant at a p-value of 0.0003146.

So we now know there is a positive relationship between cholesterol and age. However the t-test has limitations. What is the magnitude of this relationship during aging? Does it change by approximately the same amount per year? What if we don't want to break our data into groups?

How we Evaluate which Model to Use

There are a ton of models (or families of models) out there for different statistical purposes and with different assumptions. These assumptions, if violated, will give incorrect predictions. However, we might not know if these assumptions are true when selecting our model. Today we are hanging out in the top left corner, and we are going to learn the assumptions of *linear models* in general, the specific models we will be using today, and an example of each. We will trouble-shoot when assumptions fail later in the lesson.

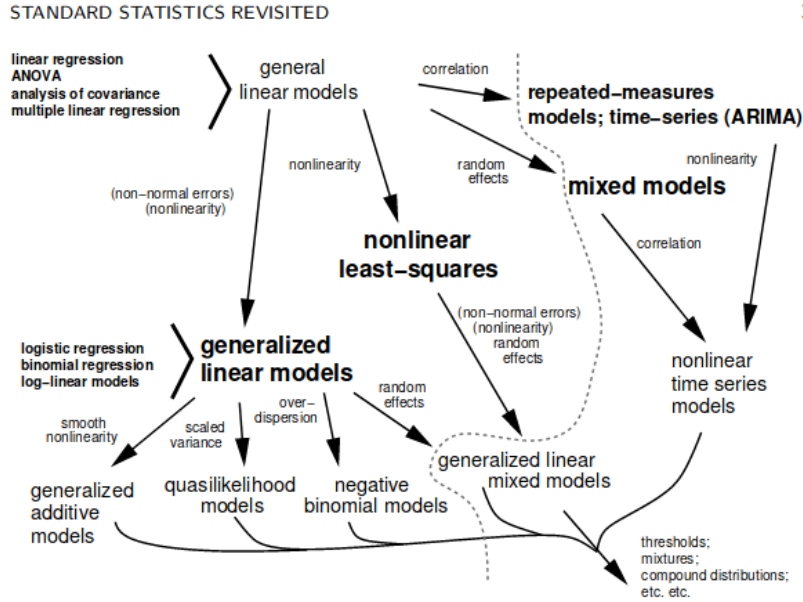


Figure 3: All (or most) of statistics. Bolker, 2007.

Distribution of Y at different x values:

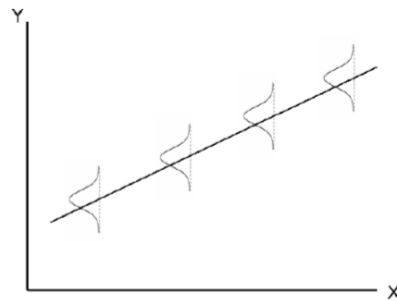


Figure 4: SISG_2016_2

Assumptions of general linear models

1. observed values are independent of each other (*independence*)
 - The probability of an event occurring does not affect the probability of another event occurring.
2. variation around expected values (residuals) are normally distributed (*normality*)
3. constant variance, homoscedastic (*equal variance*)
 - For values of x, values of y show equal variance.
4. observed values (y) are related by linear functions of the parameters to x (*linearity*)
 - Example: For $y = a + b_1x + b_2x^2$, the parameters a, b_1 , and b_2 are linear even though the independent variable has a quadratic component, x^2 . However, $y = ax^b$ is nonlinear with respect to the parameter b,

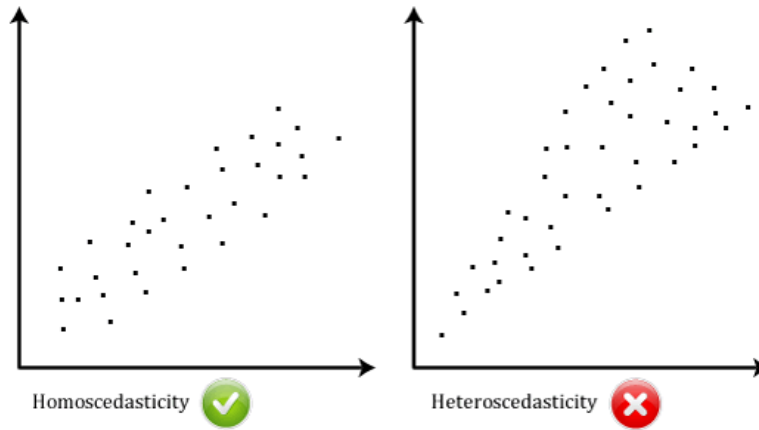


Figure 5:

and is not suitable for linear regression.

Assumptions 2 and 3 are often grouped together.

The Linear Models

For **simple linear regression** we are modelling a continuous outcome by a single continuous variable. Example: modelling cholesterol using BMI.

For **multiple linear regression** we are modelling a continuous outcome by more than one continuous variable. Example: modelling cholesterol using BMI AND age. In this case, we must consider whether there is an *interaction* between age and BMI on cholesterol (more on interactions to follow).

For **one-way ANOVA** we are modelling a continuous outcome by a single categorical variable. Example: modelling cholesterol by sex. It is important that categorical variables are explicitly input as factors to be interpreted properly in the model. For example, since we have encoded sex as 0 and 1 (instead of 'M' and 'F'), we need to specify that sex is to be treated as a categorical variable and not a number. Therefore we specify sex as a factor of 2 levels, 0 and 1.

For **multi-way ANOVA** we are modelling a continuous outcome by more than one categorical variable. Example: modelling sex and APOE genetic variants. Again, we need to consider any interaction between our categorical variables, and we need to specify our numeric values to be treated as categorical variables and not numbers. APOE will be a factor of 6 levels, one for each genetic variant.

Lastly, for **ANCOVA** we are modelling a continuous variable by a combination of categorical AND continuous variables. This could be modelling cholesterol using the genetic variants of APOE and BMI. Again, our categorical variable must be input as a factor. ANCOVA allows for each group (each genetic variant of APOE in this example) to have a separate slope.

This is a summary table you might find helpful for choosing a model based on the data types you have and the assumptions you are making. I hope to show that model selection is akin to going through mental checklist for your data, and not that scary. The independence assumption is required for all the models below, and is not included in the chart for spacing reasons.

model	categorical	continuous	linearity	normality	equal_variance
simple linear regression	X	✓	✓	✓	✓
multiple linear regression	X	✓ ✓	✓	✓	✓
one-way analysis of variance (ANOVA)	✓	X	✓	✓	✓

model	categorical	continuous	linearity	normality	equal_variance
multi-way analysis of variance (ANOVA)	✓ ✓	X	✓	✓	✓
analysis of covariance (ANCOVA)	✓	✓	✓	✓	✓
nonlinear least squares	X	✓	X	✓	✓
nonlinear analysis of covariance (ANCOVA)	✓	✓	X	✓	✓
generalized linear models	✓	✓	X*	X	X

*restricted cases

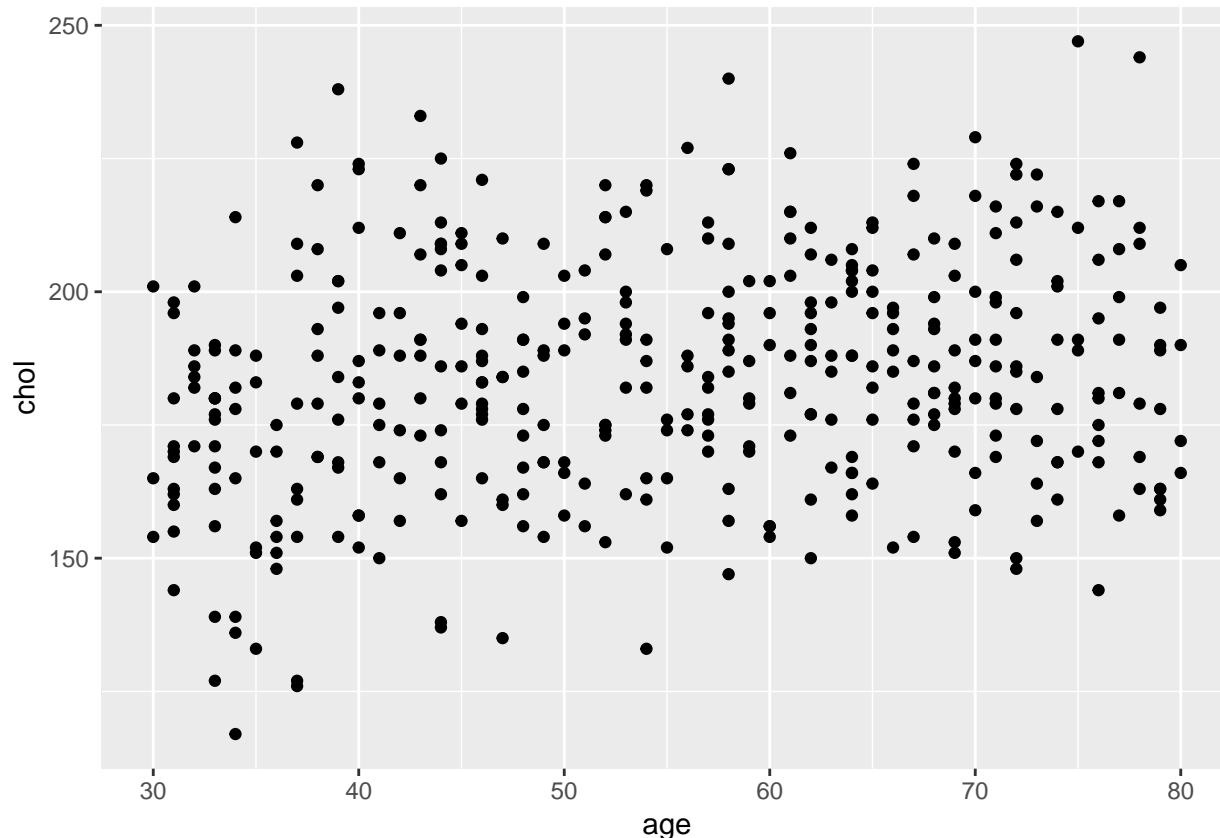
Revisiting our question:

What is the relationship between cholesterol and age?

Now, we can pick a model to answer our question instead of a t-test by considering the assumptions above.

If we evaluate our independent and dependent variables, age and cholesterol, they are both continuous, not categorical. We only have one independent variable. From the plot we made earlier (repeated here) it looks like if there is a relationship between age and cholesterol it would be linear. Data points have an even spread so the variance is likely equal and normally distributed. The values are independent (from separate blood draws).

```
ggplot(cholesterol, aes(age, chol)) +  
  geom_point()
```

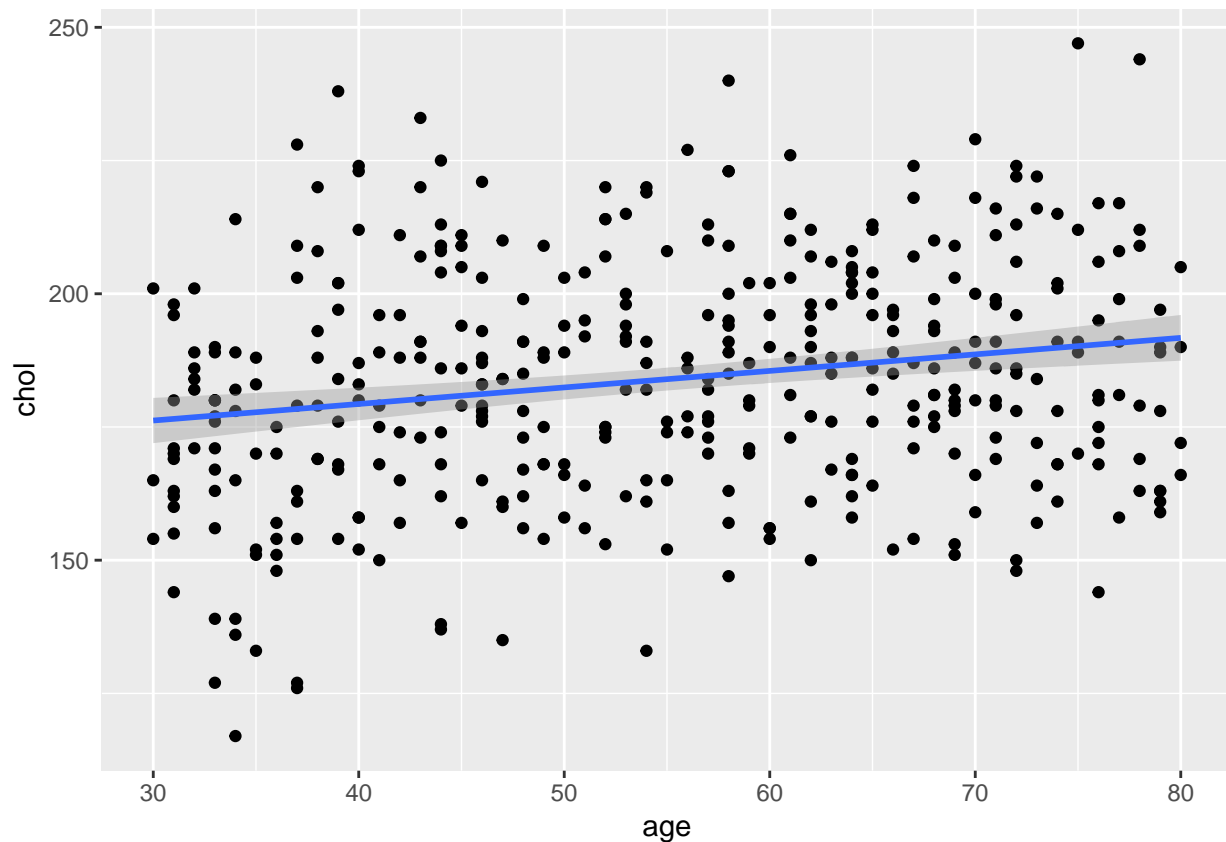


Based on the above criteria, we will try using a simple linear regression to test the association of mean serum cholesterol with age.

Simple linear regression

What we are looking for then, is the slope of the line relating cholesterol to age, which will tell us the magnitude and direction of the relationship between these variables. We can look at the slope for the linear model that `ggplot` would fit for us for an idea of what our model will look like.

```
ggplot(cholesterol, aes(age, chol)) +  
  geom_point() +  
  stat_smooth(method = "lm")
```



Review: the equation for a straight line.

Expression:

$$Y \sim \text{Normal}(a + bx, \sigma^2)$$

Y is our dependent variable that we are attempting to model. x is our independent variable.

a is the intercept (the value of y where $x = 0$; where x crosses the y -axis).

b is the slope of the line (the change in y corresponding to a unit increase in x).

Normal is telling us that our error is normally distributed.

σ^2 is the variance (squared deviation of the variable x from its mean)

Slopes

- A flat line ($b = 0$) would mean that there is no association between x and y .

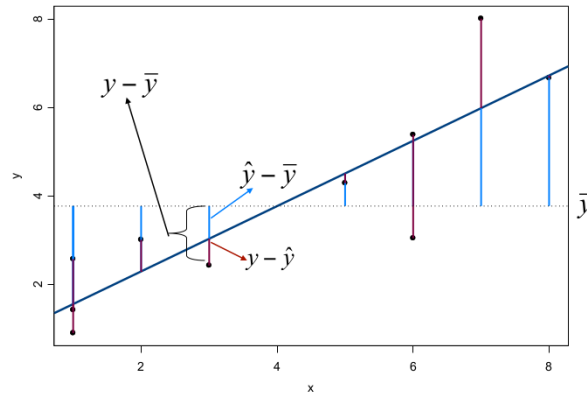


Figure 6:

- The above example has a positive slope, meaning that y increases as values of x increase.
- With a negative slope, y decreases as values of x increase.

The **interpretation** in our example is that the slope is the difference in mean serum cholesterol associated with a one year increase in age.

With a straight line we are not, of course, plotting through all of our points, but rather close to the mean of an outcome in y as a function of x . For example, there are values of cholesterol for about six 50 year-olds, and our line will fall somewhere close to the mean of these values. Values of y have a distribution at a given x , which we have assumed is normally distributed.

Lastly, in this equation we also have some normally distributed error - sampling error exists in our estimates, because different estimates give different means.

Okay, but how do we actually find the best fitting line? We use **least squares estimation**, which minimizes the sum of squares of the vertical distances from the observed points to the least squares regression line ($y - \hat{y}$).

y - observed value
 \hat{y} - estimated value
 \bar{y} - sample mean

Let's run this simple linear regression. Using R, the intercept and slope terms are implicit.

R code:

```
lm(y ~ x)
```

There are times when the intercept, the value of y at $x = 0$, doesn't make much intuitive sense to interpret.

As we are used to with writing equations, our dependent variable (cholesterol) is on the left side the `lm` formula and our independent variable (age) is on the right side; tilde `~` separates these sides. We also input the dataset to the `lm` function.

```
lm(chol ~ age, data = cholesterol)

##
## Call:
## lm(formula = chol ~ age, data = cholesterol)
##
## Coefficients:
## (Intercept)      age
##    166.9017    0.3103
```

The function will output our formula, the slope and the intercept. However, if we save the output of the function into an object, 'fit', we get a list object of the model, the input, and all associated statistics. We can look at a summary and get residuals, errors, p-values and more in addition to our coefficients.

```
fit <- lm(chol ~ age, data = cholesterol)

summary(fit)

##
## Call:
## lm(formula = chol ~ age, data = cholesterol)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -60.453 -14.643  -0.022  14.659  58.995
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 166.90168    4.26488   39.134 < 2e-16 ***
## age          0.31033     0.07524    4.125 4.52e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.69 on 398 degrees of freedom
## Multiple R-squared:  0.04099,    Adjusted R-squared:  0.03858
## F-statistic: 17.01 on 1 and 398 DF,  p-value: 4.522e-05
```

Interpretation

The intercept is 166.9 and the slope is 0.31. What does that actually mean? It means a baby (age 0) would be expected to have on average a serum cholesterol of 166.9 mg/dl. For every yearly increase in age, mean serum cholesterol is expected to increase by 0.31 mg/dl. These results are significant with a p-value < 0.001. We can reject the null hypothesis and say that mean serum cholesterol is significantly higher in older individuals. The Multiple R-squared value tells us that about 4% of the variability in cholesterol is explained by age.

We can further get confidence intervals for these values to say that 95% of the time we expect the cholesterol of a baby to fall within 158.5-175.3 mg/dl, or that we are 95% confident that the difference in mean cholesterol associated with a one year increase in age is between 0.16 and 0.46 mg/dl.

```
confint(fit)

##              2.5 %      97.5 %
## (Intercept) 158.5171656 175.2861949
## age          0.1624211   0.4582481
```

Multiple linear regression

In multiple linear regression we use multiple continuous dependent variables to predict outcome values. Additional terms can be added in 2 ways.

Adding powers of a variable (polynomial regression)

It was mentioned before that the 'linear' part of linear regression is the linear function of the *parameters* and not the independent variables. In the example below, the parameters a , b_1 , and b_2 are linear even though we

have the independent variable has a quadratic component, x^2 . An example of this could be synthesizing a chemical, where with increasing temperature synthesis progresses with an increasing curve.

Expression:

$$Y \sim \text{Normal}(a + b_1x + b_2x^2, \sigma^2)$$

If we were to write this in R, again our intercept and coefficients are implicit. To write the quadratic term we use the function `I` which just means ‘asis’.

R code:

```
lm(y ~ x + I(x^2))
```

Adding extra variables to our model

We are interested in improving our model by adding extra variable we think might have an effect on our outcome values. In the example below, we are adding the independent variables x_1 , x_2 , and each of these terms has their own linear parameter b_1 and b_2 , respectively.

Expression:

$$Y \sim \text{Normal}(a + b_1x_1 + b_2x_2, \sigma^2)$$

This is the model we will be using next. To aid with interpretation let’s think about a parameter. b_2 is the expected mean change in unit per change in x_2 if x_1 is held constant (sometimes called controlling for x_1).

The null hypothesis in this case is that all $b_1, b_2 = 0$. The alternative hypothesis is that at least one of these parameters is not null.

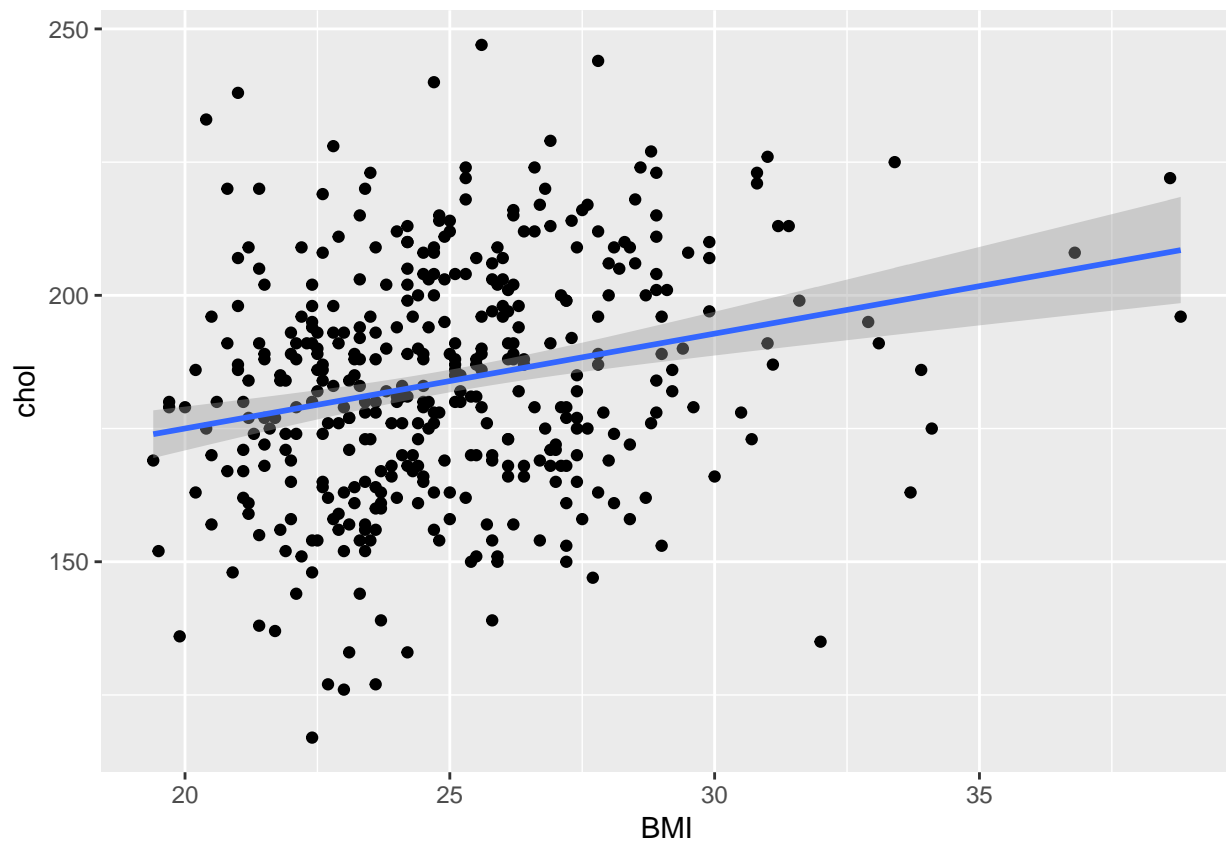
Again in R the intercept and coefficients are implicit in the the `lm` function.

R code:

```
lm(y ~ x1 + x2)
```

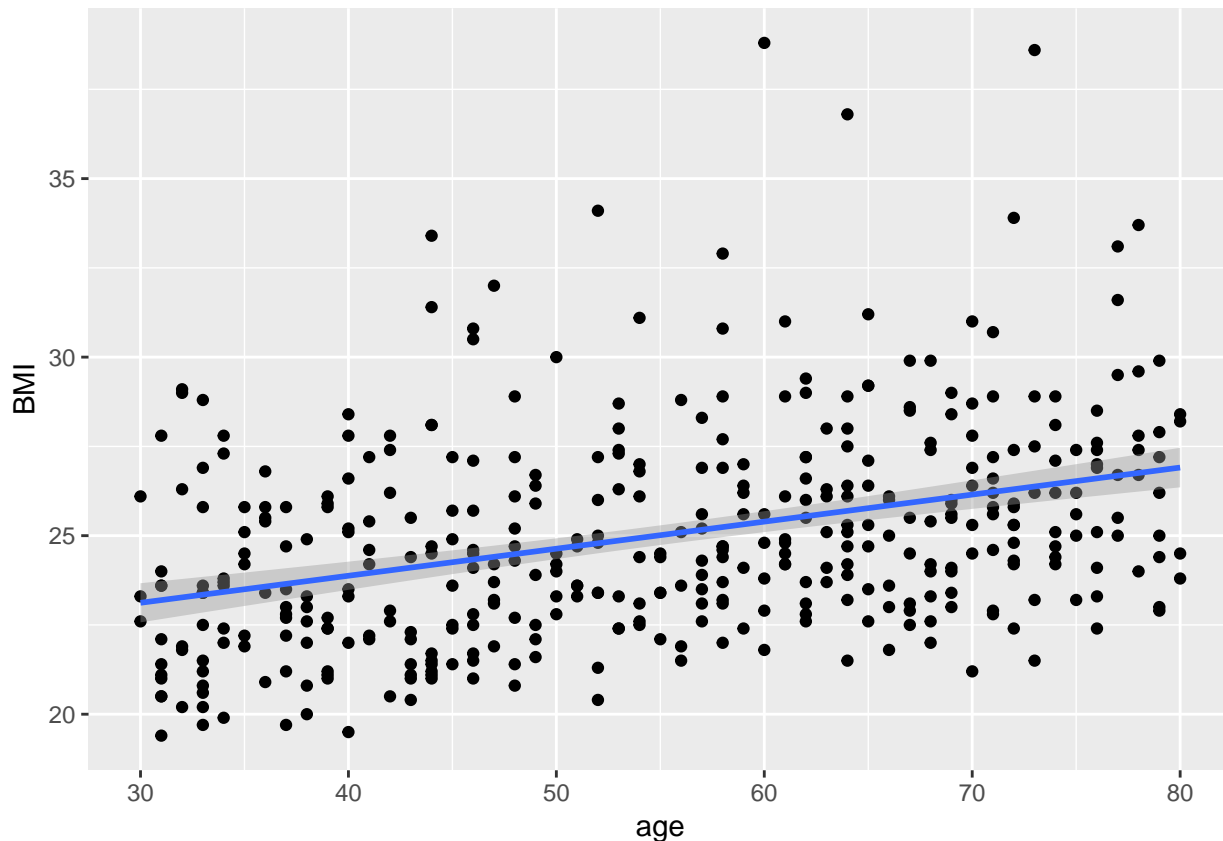
We know that age has an effect on cholesterol. With our new model we want to ask the question: **Is there a statistically significant relationship between mean serum cholesterol and age after controlling for BMI?** Let’s look graphically at these relationships to help us understand our model. First let’s plot BMI vs cholesterol. We can add a linear fit to make sure we are expecting a positive slope.

```
ggplot(cholesterol, aes(BMI, chol)) +  
  geom_point() +  
  stat_smooth(method = "lm")
```



We should also take a look at the relationship between BMI and age.

```
ggplot(cholesterol, aes(age, BMI)) +  
  geom_point() +  
  stat_smooth(method = "lm")
```

Cholesterol increases with BMI. BMI increases with age. We will look at the association of age and cholesterol while holding BMI constant to see if the significance of our finding of the increase in cholesterol with age was affected by BMI.

```
mfit <- lm(chol ~ age + BMI, data = cholesterol)
```

```
summary(mfit)
```

```
##
## Call:
## lm(formula = chol ~ age + BMI, data = cholesterol)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -58.994 -15.793   0.571  14.159  62.992
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  137.1612    9.0061  15.230 < 2e-16 ***
## age           0.2023     0.0795   2.544 0.011327 *
## BMI           1.4266     0.3822   3.732 0.000217 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.34 on 397 degrees of freedom
## Multiple R-squared:  0.07351,    Adjusted R-squared:  0.06884
## F-statistic: 15.75 on 2 and 397 DF,  p-value: 2.62e-07
```

Interpretation

Our equation would now look like $y = 137.16 + 0.20age + 1.43BMI$.

The estimated increase in mean serum cholesterol over after one year holding BMI constant is 0.20 mg/dl. This increase is less than our previous value of 0.31 mg/dl. Why do the estimates differ?

Before, we were not controlling for BMI. Our estimates of the age associated increase in mean cholesterol is now for subjects with the *same* BMI and not for subjects with *all* BMIs.

It looks like both age and BMI are significant. But we might want to verify - did adding BMI actually make a difference to the model?

We can compare these models with the `anova` function. The output of our model, 'mfit', is an `lm` object. With 2 `lm` objects, the `anova` function tests the *models* against one another to see if their coefficients are significantly different and prints these results in an analysis of variance table. (Given 1 `lm` object, it will test whether model *terms* of a model are significant - we will be using the function in this format later.)

```
anova(fit, mfit)
```

```
## Analysis of Variance Table
##
## Model 1: chol ~ age
## Model 2: chol ~ age + BMI
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      398 187187
## 2      397 180842  1    6345.8 13.931 0.0002174 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Our second model is a significantly different from our first model. What is this significance based on?

The significance is a probability based on an F-test. While the t-test tells if you a single variable is statistically significant, and an **F-test** tells you if a group of variables is jointly significant. Since the F-test compares the joint effect of all variables together, a large F value means 'something' is significant. F-tests are not used alone because you still need to use a p-value to find out 'what' is significant.

Interaction terms

What is meant by an interaction? There is an **interaction** if the association between the response and the predictor variable changes across the range of the new variable. This can be seen in the expression below, where the difference in means between x_1 and x_2 changes additionally by b_3 for each unit difference in x_2 or x_1 , ie. the slope of x_1 changes with x_2 , because b_3 is changing.

Expression:

$$Y \sim \text{Normal}(a + b_1x_1 + b_2x_2 + b_3x_1x_2, \sigma^2)$$

In the graph below, there is an interaction between education and ideology. The slope indicating the probability that people will care if sea level rises 20 feet, changes with each education level and each shift in ideology. If there was no interaction with ideology and education, the slopes shown would be parallel.

When testing for an interaction between 2 input variables, the `lm` input takes an asterik '*' instead of a plus sign between the dependent variables.

R code:

```
lm(y ~ x1*x2)
```

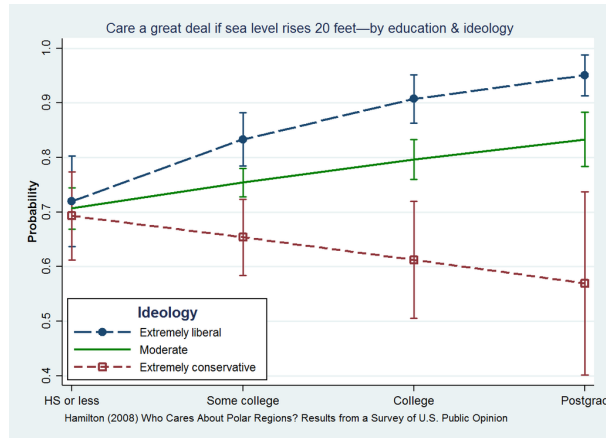


Figure 7:

An interaction is different than a **confounding factor**, for which the association between the response and predictor variable is constant across the range of the new variable. You can think of them as variables that have an effect on the outcome, but haven't been accounted for. For example, in our first model where the increase in cholesterol was ONLY due to an increase in age, BMI would be a confounding factor because weight contributes significantly to an increase in cholesterol, and age alone is not responsible for the increase in cholesterol.

Challenge

Test if there is an interaction between age and BMI in a model predicting mean serum cholesterol. Is the interaction significant? Is there a difference between this model and the model with age as the only variable? Is there a difference between this model and the model of BMI and age model with no interaction?

One-way analysis of variance (ANOVA)

In the analysis of variance (ANOVA) independent variables are categorical (factors) rather than continuous. This allows us to ask the the question:

Does the genetic factor rs174548 have an effect on cholesterol levels?

Our categorical example is represented by α_i . i represents the levels of our factor.

Expression:

$$Y \sim \text{Normal}(\alpha_i, \sigma^2)$$

We still use the `lm` function, however we replace our continuous variable with `f`, a categorical variable (factor). If your data is character type, R will automatically make a factor for you. However if your data is numeric, R will interpret it as continuous. In this case, you need to make your numeric data a factor first using `factor`.

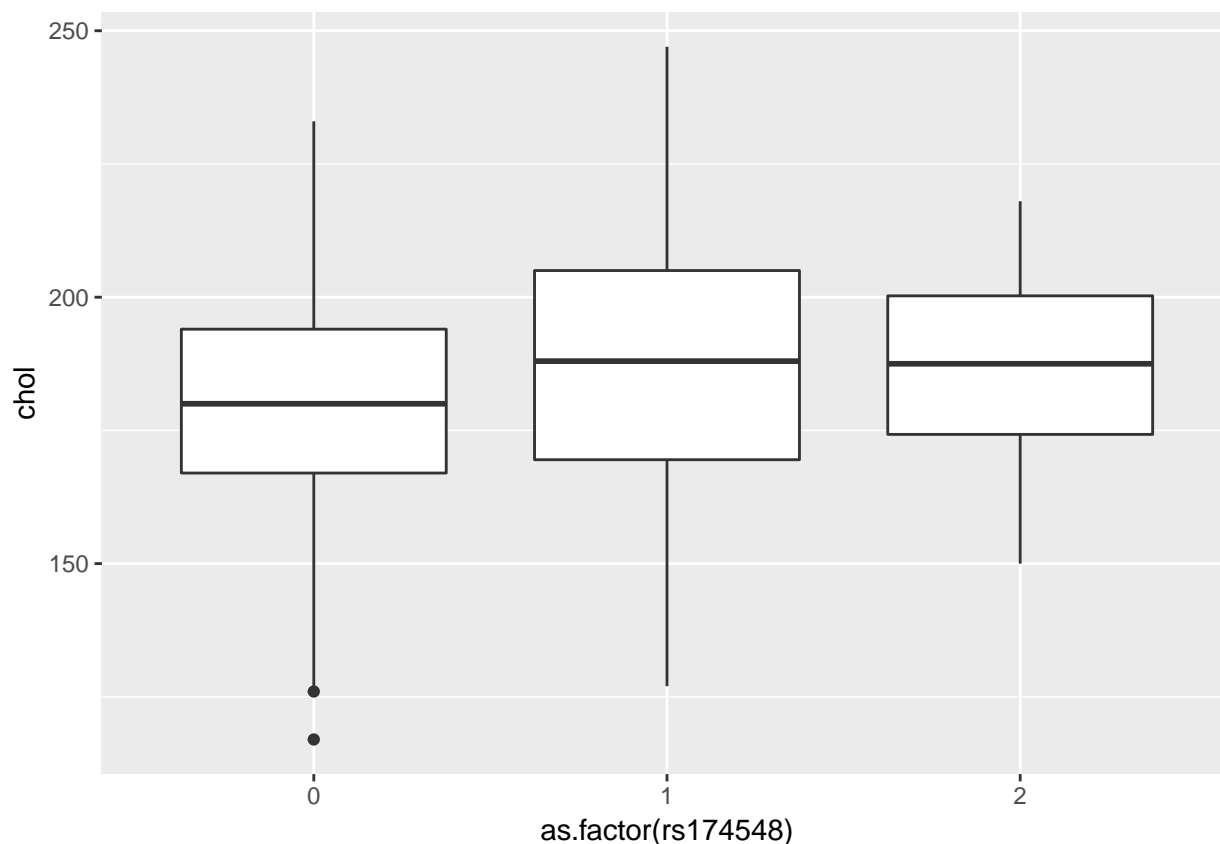
R code:

```
lm(y ~ f)
```

R parameterizes the model in terms of the differences between the first group and subsequent groups (ie. relative to the first group) rather than in terms of the mean of each group. This is similar to the interpretation of the previous linear models. (You can instead fit the means of each group using: `lm(y ~ f-1)`).

To begin to answer our question, we can first plot the relationship between rs174548 and cholesterol.

```
ggplot(cholesterol, aes(as.factor(rs174548), chol)) + geom_boxplot()
```



Our genetic factor has 3 groups, and we will be comparing the means for each of these groups. These groups have high variance, and there is a good deal of overlap between them.

To assess whether the means are equal, the model compares:

- variation between the sample means (MSR)
- natural variation of the observations within the sample (MSE)

The larger the MSR compared to the MSE the more support there is for a difference between the population means. The ratio of MSR/MSE is our F-statistic.

Dummy Variables

We can encode our categorical variable as a *dummy variable*. 0 in our data frame stands for the genotype C/C, 1 is C/G and 2 is G/G. But instead we can create a matrix with $k - 1$ separate columns of 0's and 1's, where k is the number of factor levels in our categorical variable. The omitted category is the reference group. Each genetic factor has a unique encoding where 0 means the SNP is not present and 1 means that SNP is present.

rs174548	x1	x2
C/C	0	0

rs174548	x1	x2
C/G	1	0
G/G	0	1

Now we can do regression with our dummy variables. The form of this equation should look more familiar.

Expression:

$$Y \sim \text{Normal}(\beta_0 + \beta_1 x_1 + \beta_2 x_2, \sigma^2)$$

Interpretation

The interpretation of this model is a bit trickier.

- β_0 - mean cholesterol when rs174548 is C/C
- $\beta_0 + \beta_1$ - mean cholesterol when rs174548 is C/G
- $\beta_0 + \beta_2$ - mean cholesterol when rs174548 is G/G

Alternatively,

- β_1 is the difference in mean cholesterol levels between groups with rs174548 equal to C/G and C/C
- β_2 is the difference in mean cholesterol levels between groups with rs174548 equal to G/G and C/C

So you can think of each of these groups having their own means. ie. $\mu_0 = \beta_0$, $\mu_1 = \beta_0 + \beta_1$, $\mu_2 = \beta_0 + \beta_2$. We are testing the hypothesis whether these means are equal or not.

R will create dummy variables in the background if you state you have a categorical variable.

```
anfit1 <- lm(chol ~ as.factor(rs174548), data = cholesterol)
summary(anfit1)
```

```
##
## Call:
## lm(formula = chol ~ as.factor(rs174548), data = cholesterol)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -64.062 -15.913  -0.062  14.938  59.136
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      181.062      1.455  124.411  < 2e-16 ***
## as.factor(rs174548)1    6.802      2.321   2.930  0.00358 **
## as.factor(rs174548)2    5.438      4.540   1.198  0.23167
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.93 on 397 degrees of freedom
## Multiple R-squared:  0.0221, Adjusted R-squared:  0.01718
## F-statistic: 4.487 on 2 and 397 DF,  p-value: 0.01184
```

Interpretation

- The intercept, 181.06 mg/dl is the mean cholesterol when rs174548 is C/C.
- 181.06 mg/dl + 6.80 mg/dl is the mean cholesterol when rs174548 is C/G.
- 181.06 mg/dl + 5.44 mg/dl is the mean cholesterol when rs174548 is G/G.

Alternatively,

- 6.80 mg/dl is the difference in mean cholesterol levels between groups with rs174548 equal to C/G and C/C
- 5.44 mg/dl is the difference in mean cholesterol levels between groups with rs174548 equal to G/G and C/C

The genetic factor rs174548 has an effect on cholesterol at a significance level of $p < 0.05$.

An analysis of variance table with one model as input test gives us the same p-value. Why?

```
anova(anfit1)

## Analysis of Variance Table
##
## Response: chol
##              Df Sum Sq Mean Sq F value    Pr(>F)
## as.factor(rs174548)    2    4314   2157.10    4.4865 0.01184 *
## Residuals              397   190875    480.79
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This tells us that there is a difference in means (rejects the null hypothesis that all means are the same), but doesn't tell us which means are different.

In order to look at this we need to look at multiple pairwise comparisons.

$$\mu_0 = \mu_1, \mu_0 = \mu_2, \mu_1 = \mu_2$$

Multiple test correction

Multiple comparisons increase the **family-wise error rate (FWER)** - the probability of making a false discovery (aka a false positive or Type I error). This is where multiple test corrections come in to control the error at a specific threshold (ie. $\alpha = 0.05$ or 5%). One of the simplest and conservative is the Bonferroni correction (α/k or multiplying p-values by k).

Previously, it was mentioned that you can fit the means for each group using $\text{lm}(y \sim f-1)$. To run multiple tests to see if group means differ we can use this equation for **general linear hypothesis testing**, which takes in a model as well as a **contrast matrix** for the comparisons you want to make. The simplest contrast matrix is a matrix of 0, 1, and -1's where the relationship -1 and 1 are the factor levels for which you want to test the differences.

```
tfit <- lm(chol ~ -1 + as.factor(rs174548), data = cholesterol)

M <- contrMat(table(cholesterol$rs174548), type = "Tukey")
M
```

```
##
## Multiple Comparisons of Means: Tukey Contrasts
##
##      0  1  2
## 1 - 0 -1  1  0
## 2 - 0 -1  0  1
## 2 - 1  0 -1  1
```

This rownames of the contrast matrix tell us what is being compared ([1-0], [2-0], [2-1]). For example [1-0] is the difference between C/C (-1) and C/G (1). More complicated comparisons can be made. For example, the difference between C/C and the average of G/G and C/G could be specified by adding a row to the matrix of -2 1 1 (Note: rows of a contrast matrix must add to zero).

To get estimates using general linear hypothesis testing we use the `glht` function; our linear hypotheses to be tested are specified by our contrast matrix. We will first look at a summary without adjusting/correcting our p-values.

```
mc <- glht(tfit, linfct = M)

summary(mc, test = adjusted("none"))

##
## Simultaneous Tests for General Linear Hypotheses
##
## Multiple Comparisons of Means: Tukey Contrasts
##
##
## Fit: lm(formula = chol ~ -1 + as.factor(rs174548), data = cholesterol)
##
## Linear Hypotheses:
##           Estimate Std. Error t value Pr(>|t|)
## 1 - 0 == 0    6.802      2.321   2.930 0.00358 **
## 2 - 0 == 0    5.438      4.540   1.198 0.23167
## 2 - 1 == 0   -1.364      4.665  -0.292 0.77015
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## (Adjusted p values reported -- none method)
```

Interpretation

[1 - 0 == 0] The difference in means between C/C and C/G is 6.80 mg/dl and this difference is significant.
 [2 - 0 == 0] The difference in means between C/C and G/G is 5.44 mg/dl and this difference is not significant.
 [2 - 1 == 0] The difference in means between C/G and G/G is -1.36 mg/dl and this difference is not significant.

We can see if multiple test correction affects these relationships.

```
summary(mc, test = adjusted("bonferroni"))

##
## Simultaneous Tests for General Linear Hypotheses
##
## Multiple Comparisons of Means: Tukey Contrasts
##
##
## Fit: lm(formula = chol ~ -1 + as.factor(rs174548), data = cholesterol)
##
## Linear Hypotheses:
##           Estimate Std. Error t value Pr(>|t|)
## 1 - 0 == 0    6.802      2.321   2.930 0.0107 *
## 2 - 0 == 0    5.438      4.540   1.198 0.6950
## 2 - 1 == 0   -1.364      4.665  -0.292 1.0000
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## (Adjusted p values reported -- bonferroni method)
## #other correction types
## summary(mc, test = adjusted("BH"))
## summary(mc, test = adjusted("fdr"))
```

The significant difference in mean cholesterol between C/C and C/G genotypes of rs174548 holds under different multiple test corrections.

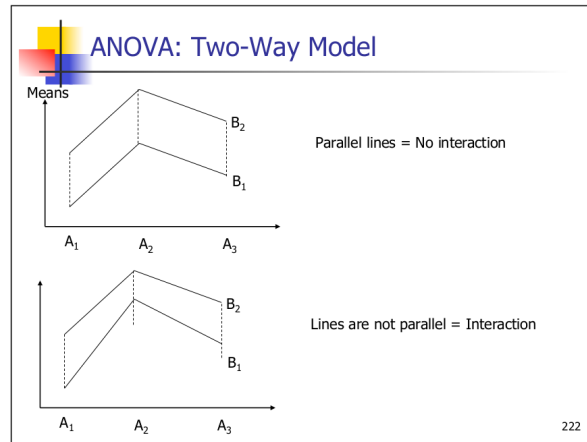


Figure 8: SISG_2016

Multi-way analysis of variance (ANOVA)

Two or more categorical variables (factors) are used to model our outcome. We can now ask the question:

Does the effect of the genetic factor rs174548 differ between males and females?

We need to test whether there is an effect of our factors on cholesterol and also if there is an interaction between these factors.

Expression:

$$Y \sim \text{Normal}(\alpha_i + \beta_j, \sigma^2)$$

α and β are our categorical variables. i is the level of the first group, and j is the level of the second group.

As with one-way ANOVA, R models our categorical variables as factors.

R code:

`lm(y ~ f1 + f2)`, testing for main effects without interaction.

`lm(y ~ f1*f2)`, testing for the main effects with interaction.

The following diagram will help us visualize the differences in coefficients with and without interaction between 2 categorical variables.

In this first scenario, the difference in the means between groups defined by factor B does not depend on the level of factor A and vice versa. This means that there is no interaction, and the lines between the factor groups are parallel. In the second scenario the difference in the means between groups defined by factor B changes when A₂ is present. There is an interaction and the lines are not parallel.

We can first run a two-way model without testing for interaction.

```
twofit <- lm(chol ~ as.factor(sex) + as.factor(rs174548), data = cholesterol)
summary(twofit)
```

```
##
## Call:
## lm(formula = chol ~ as.factor(sex) + as.factor(rs174548), data = cholesterol)
##
## Residuals:
```



```
##      Min      1Q  Median      3Q      Max
## -66.653 -14.463  -0.601  15.445  57.635
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      175.365      1.786  98.208 < 2e-16 ***
## as.factor(sex)1       11.053      2.126   5.199 3.22e-07 ***
## as.factor(rs174548)1    7.236      2.250   3.215 0.00141 **
## as.factor(rs174548)2    5.184      4.398   1.179 0.23928
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.24 on 396 degrees of freedom
## Multiple R-squared:  0.08458,    Adjusted R-squared:  0.07764
## F-statistic: 12.2 on 3 and 396 DF,  p-value: 1.196e-07
```

Interpretation

- The estimated mean cholesterol for males in C/C group is the intercept, 175.36 mg/dl.
- The estimated difference in mean cholesterol between females and males controlled for genotype is 11.05 mg/dl.
- The estimated difference in mean between C/G and C/C groups controlled for gender is 7.24 mg/dl.
- The estimated difference in mean between G/G and C/C groups controlled for gender is 5.18 mg/dl.

There is evidence cholesterol is associated with gender ($p < 0.001$).

How does this compare to the model with gender alone as a predictor?

```
genfit <- lm(chol ~ as.factor(sex), data = cholesterol)
anova(genfit, twofit)
```

```
## Analysis of Variance Table
##
## Model 1: chol ~ as.factor(sex)
## Model 2: chol ~ as.factor(sex) + as.factor(rs174548)
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      398 183480
## 2      396 178681  2    4799.1 5.318 0.005259 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

There is a difference between these 2 models ($p = 0.005$).

We can now check the two-way anova with the interaction.

```
intfit2 <- lm(chol ~ as.factor(sex) * as.factor(rs174548), data = cholesterol)
summary(intfit2)
```

```
##
## Call:
## lm(formula = chol ~ as.factor(sex) * as.factor(rs174548), data = cholesterol)
##
## Residuals:
##      Min      1Q  Median      3Q      Max
## -70.529 -13.604  -0.974  14.171  54.882
##
```

```
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   178.1182     2.0089  88.666 < 2e-16
## as.factor(sex)1                 5.7109     2.7982   2.041  0.04192
## as.factor(rs174548)1            0.9597     3.1306   0.307  0.75933
## as.factor(rs174548)2           -0.2015     6.4053  -0.031  0.97492
## as.factor(sex)1:as.factor(rs174548)1 12.7398     4.4650   2.853  0.00456
## as.factor(sex)1:as.factor(rs174548)2 10.2296     8.7482   1.169  0.24297
##
## (Intercept)                    ***
## as.factor(sex)1                 *
## as.factor(rs174548)1
## as.factor(rs174548)2
## as.factor(sex)1:as.factor(rs174548)1 **
## as.factor(sex)1:as.factor(rs174548)2
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.07 on 394 degrees of freedom
## Multiple R-squared:  0.1039, Adjusted R-squared:  0.09257
## F-statistic:  9.14 on 5 and 394 DF,  p-value: 3.062e-08
```

Interpretation

- The estimated mean cholesterol for males in the C/C group is 178.12 mg/dl.
- The estimated mean cholesterol for females in the C/C group is $(178.12 + 5.71)$ mg/dl.
- The estimated mean cholesterol for men in the C/G group $(178.12 + 0.96)$ mg/dl.
- The estimated mean cholesterol for females in the C/G group is $(178.12 + 5.71 + 0.96 + 12.74)$ mg/dl.

There appears to be a significant interaction between being female and having the C/G genotype.

Let's compare the with interaction and without interaction model.

```
anova(twofit,intfit2)
```

```
## Analysis of Variance Table
##
## Model 1: chol ~ as.factor(sex) + as.factor(rs174548)
## Model 2: chol ~ as.factor(sex) * as.factor(rs174548)
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      396 178681
## 2      394 174902  2    3778.9 4.2564 0.01483 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

There is evidence that these two models are different ($p = 0.015$).

Analysis of covariance (ANCOVA)

The analysis of covariance (ANCOVA) model allows for different intercepts and slopes with respect to a continuous variable in different categorical groups. ANCOVA, therefore, has a linear regression component. This allows us to ask the question:

Is the relationship between age and cholesterol is affected by gender?

Expression:

$$Y \sim \text{Normal}(\beta_i + \beta_i x, \sigma^2)$$

Parameters are the intercept of the first factor level, the slope with respect to x for the first factor level, the differences in the intercepts for each factor level other than the first, and the differences in slopes for each factor level other than the first.

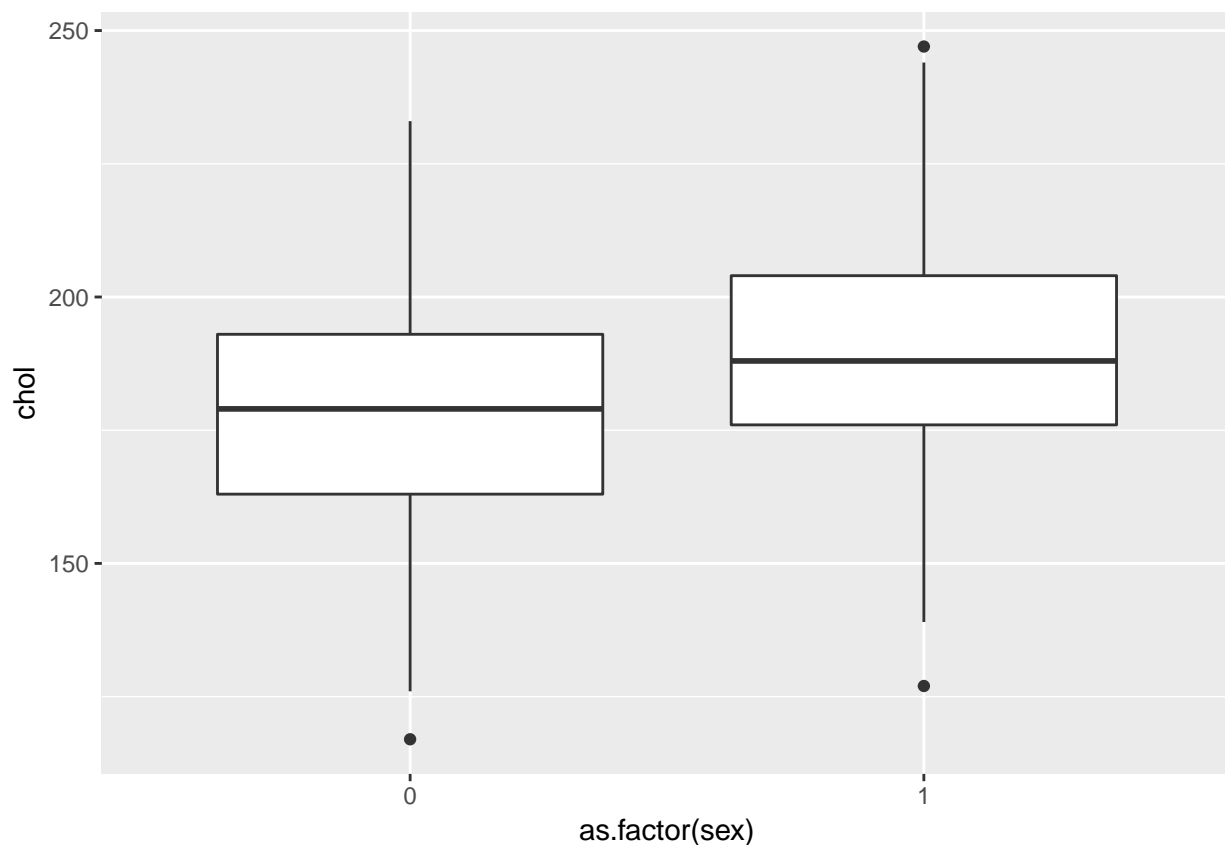
R code:

`lm(y ~ f + x)`, testing for main effects without interaction.

`lm(y ~ f*x)`, testing for main effects with interaction.

To answer our question, let's first take a quick look at gender differences in cholesterol in our dataset, keeping in mind that males are encoded as 0 and females as 1. Based on sex information alone, we see that women have a higher mean serum cholesterol, but we don't know if this is significant.

```
ggplot(cholesterol, aes(as.factor(sex), chol)) +
  geom_boxplot()
```



Our model won't look that different from our other equations except that we have categorical and continuous predictor variables.

```
mfit2 <- lm(chol ~ age + sex, data = cholesterol)
```

```
summary(mfit2)
```

```
##
## Call:
## lm(formula = chol ~ age + sex, data = cholesterol)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -55.662 -14.482  -1.411  14.682  57.876
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 162.35445    4.24184  38.275  < 2e-16 ***
## age         0.29697     0.07313   4.061 5.89e-05 ***
## sex         10.50728     2.10794   4.985 9.29e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.06 on 397 degrees of freedom
## Multiple R-squared:  0.09748,    Adjusted R-squared:  0.09293
## F-statistic: 21.44 on 2 and 397 DF,  p-value: 1.44e-09
```

Interpretation

Controlling for sex, mean cholesterol increases by 0.30 mg/dl for an additional year of age. This is close to the slope for our model of cholesterol alone, 0.31 mg/dl. This does not necessarily mean that the age/cholesterol relationship is the same in males and females; we need to check out the interaction term. There appears to be an increase in mean serum cholesterol of 10.5 mg/dl in females over males.

```
intfit <- lm(chol ~ age * sex, data = cholesterol)
```

```
summary(intfit)
```

```
##
## Call:
## lm(formula = chol ~ age * sex, data = cholesterol)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -56.474 -14.377  -1.215  14.764  58.301
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 160.31151    5.86268  27.344  < 2e-16 ***
## age         0.33460     0.10442   3.204  0.00146 **
## sex         14.56271     8.29802   1.755  0.08004 .
## age:sex      -0.07399     0.14642  -0.505  0.61361
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.08 on 396 degrees of freedom
## Multiple R-squared:  0.09806,    Adjusted R-squared:  0.09123
## F-statistic: 14.35 on 3 and 396 DF,  p-value: 6.795e-09
```

Interpretation

Males are coded as 0 and females are coded as 1 in this model. The intercept term is the mean serum cholesterol for MALES at age 0. The slope term for age is the difference in mean cholesterol associated with a one year change in age for MALES. The slope for sex is the difference in mean cholesterol between males and females at age 0. The interaction term is the difference in the change in mean cholesterol associated with each one year change in age for females compared to males. Sex exerts a small and not statistically significant effect on the age/cholesterol relationship.

Let's compare our models with and without an interaction term with anova (an F-test).

```
anova(mfit2, intfit)
```

```
## Analysis of Variance Table
##
## Model 1: chol ~ age + sex
## Model 2: chol ~ age * sex
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1     397 176162
## 2     396 176049  1    113.52 0.2554 0.6136
```

Adding the interaction term did not change the model significantly.

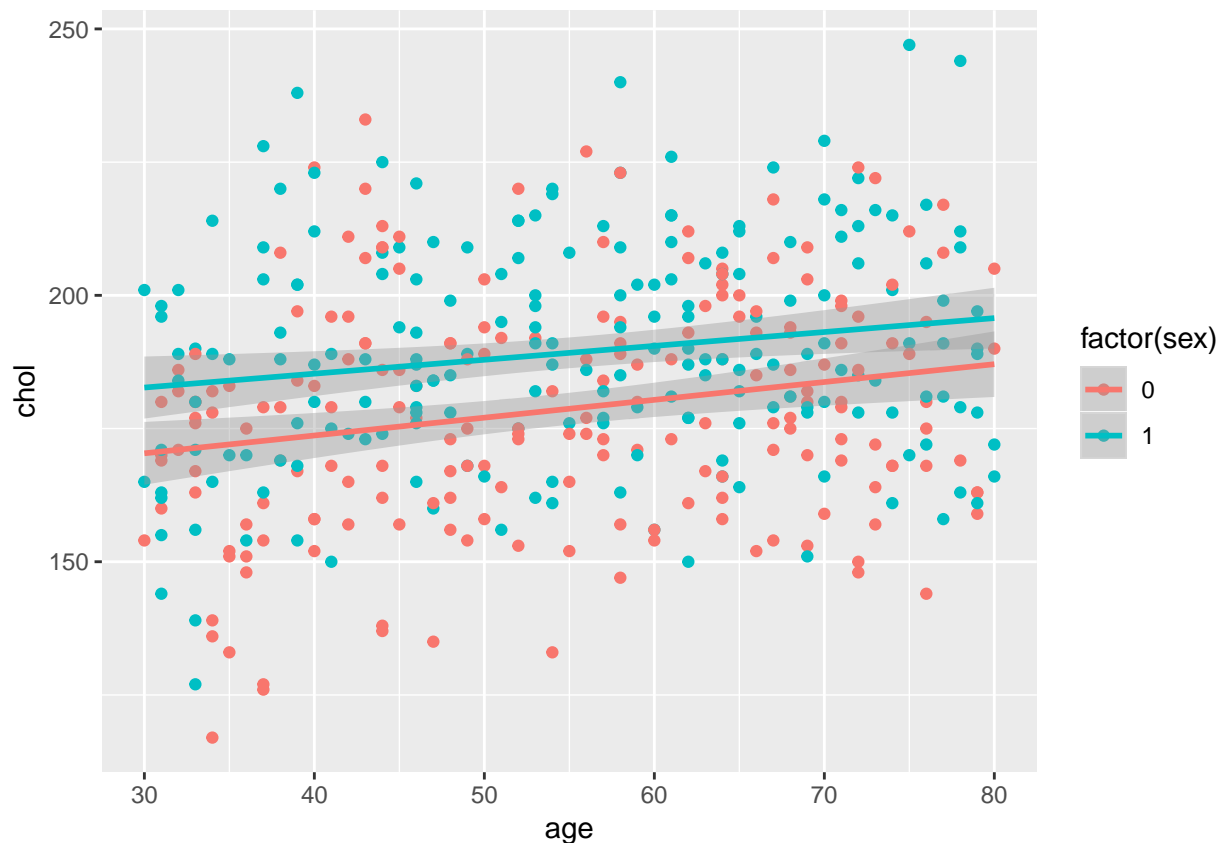
Let's compare the model with age only to the model where sex is taken into account without interaction.

```
anova(fit, mfit2)
```

```
## Analysis of Variance Table
##
## Model 1: chol ~ age
## Model 2: chol ~ age + sex
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     398 187187
## 2     397 176162  1    11025 24.846 9.291e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The addition of sex makes a statistically significant difference to the model. Take a look at the whole picture in graphical format.

```
ggplot(cholesterol, aes(age, chol, color = factor(sex))) +
  geom_point() +
  stat_smooth(method = "lm")
```



Interpretation

Gender doesn't change the relationship between age and cholesterol; these lines are almost parallel (another way to put it is whether you are male or female your cholesterol will on average be increasing by 0.30 mg/dl a year), but there is a different mean serum cholesterol estimate for males vs females that differs by 10.51 dg/ml.

Review: Models we used today

Before we move on I want to take a step back and quickly review the models and code we've gone through today. Firstly, with our example dataset, and then more generally. I hope you can see that though conceptually different, getting a handle on the code isn't too bad.

For all of these models we are trying to determine the effect of different variables on cholesterol. The differences are whether we are using continuous data, categorical data, a mixture of data types, and whether there is an interaction (*) between our input variables.

We have started with models that assume normally distributed errors, and we will investigate models with non-normal errors in a future lesson.

model	categorical	continuous	R_code
simple linear regression	X	✓	<code>lm(chol ~ age)</code>
multiple linear regression	X	✓ ✓	<code>lm(chol ~ age + BMI)</code> , <code>lm(chol ~ age*BMI)</code>
one-way analysis of variance (ANOVA)	✓	X	<code>lm(chol ~ factor(rs174548))</code>

model	categorical	continuous	R_code
multi-way analysis of variance (ANOVA)	✓ ✓	X	lm(chol ~ factor(sex) + factor(rs174548)), lm(chol ~ factor(sex)*factor(rs174548))
analysis of covariance (ANCOVA)	✓	✓	lm(chol ~ factor(sex) + age), lm(chol ~ factor(sex)*age)

In the table below, our R code for each of the models has been generalized. Here, y is our predictor variable, x is a continuous variable, and f is a categorical variable (factor).

model	R_code
simple linear regression	lm(y ~ x)
multiple linear regression	lm(y ~ x + I(x^2)), lm(y ~ x ₁ + x ₂), lm(y ~ x ₁ *x ₂)
one-way analysis of variance (ANOVA)	lm(y ~ f)
multi-way analysis of variance (ANOVA)	lm(y ~ f ₁ + f ₂), lm(y ~ f ₁ *f ₂)
analysis of covariance (ANCOVA)	lm(y ~ f + x), lm(y ~ f*x)

You need not memorize any of these charts - you may just want to use them to orient yourself in the future. Much of the R code seems the same whether you are doing multiple linear regression, ANOVA or ANCOVA, so it is good to have a reference point.

Challenge

Does the effect of the genetic factor rs174548 differ depending on a subject's age? Make a plot of age versus cholesterol and color points by genotype. Add a linear model to the plot. Are you expecting an interaction based on this plot? Test models for the association between cholesterol and age controlling for the genetic factor rs174548 with interaction and without interaction. Look at the summary statistics for each model fit. How would you interpret the results? Compare the two models with an analysis of variance table.

Challenge:

Does the genetic factor APOE have an effect on cholesterol levels? If so, does this interaction vary depending on a subject's age? Plot the relationship between APOE and cholesterol. Choose your model(s). Interpret your summary statistics. What model did you find 'best' for the job? Perform multiple hypothesis testing for the relationship between cholesterol and APOE variants. Which relationships are significant? Do these hold after multiple test correction?

Resources:

https://github.com/ttimbers/lm_and_glm/blob/master/lm_and_glm.Rmd
<https://github.com/seananderson/glmm-course>
<http://michael.hahsler.net/SMU/EMIS7331/R/regression.html>
<https://ms.mcmaster.ca/~bolker/emdbook/book.pdf>
<http://www.differencebetween.net/science/mathematics-statistics/difference-between-ancova-and-regression/>
<https://stats.stackexchange.com/questions/77563/linear-model-fit-seems-off-in-r>
<https://www.biostat.washington.edu/suminst/archives/SISG2016/SM1604>
<https://ms.mcmaster.ca/~bolker/>
<http://www.mathnstuff.com/math/spoken/here/2class/90/htest.htm>

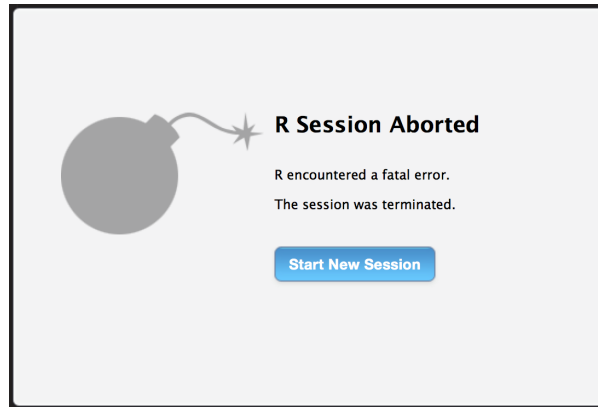


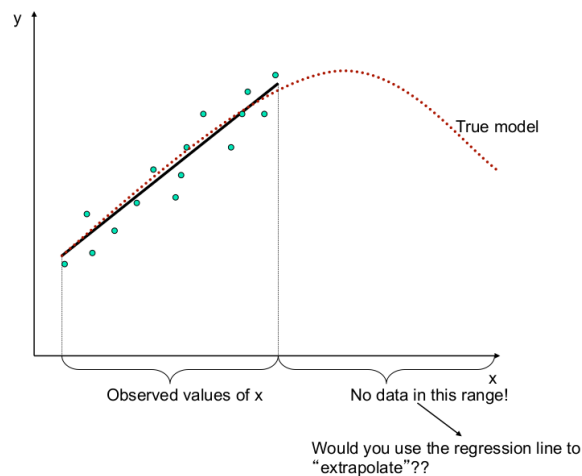
Figure 9:

Thanks for coming!!!

Appendix

Prediction

When predicting values you are assuming that your model is true. This might be fair within the range of your data. This is to be interpreted with caution outside the range of your data. For example, polynomial data may look linear over a certain range.



This is one of my favourite xkcd comics, probably just because I am getting married soon and this extrapolation is terrifying.

The `predict` function works with many different kinds of fits: not just linear models but nonlinear, polynomial, generalized linear models, etc. `predict` will try to guess the fit based on the object input, but this information can be specified using `predict.lm`. The help page for `predict.lm` is more useful as it is specific for the linear model fit.

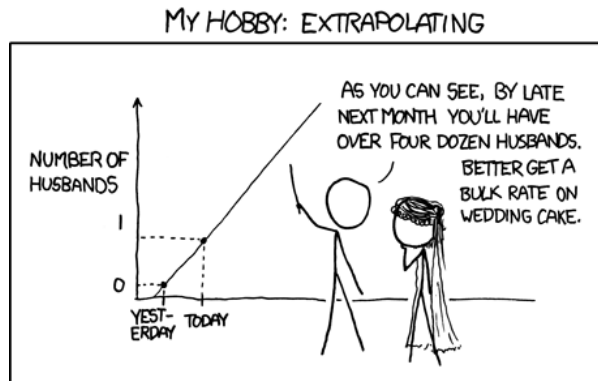


Figure 10: xkcd

Challenge

Use `predict.lm` to predict the mean cholesterol at age 47 from our model object ‘fit’. (‘fit’ is our first model, `lm(chol ~ age, data = cholesterol)`).

In addition to the linear model, the function needs the `newdata` that we want to predict. Note that `newdata` takes in a data frame. We can predict the mean cholesterol at age 47 within a confidence interval that can be specified using `level`. The output is the mean, as well as the upper and lower boundaries of the estimate.

```
predict.lm(fit, newdata = data.frame(age=47), interval = "confidence")
```

```
##          fit      lwr      upr
## 1 181.4874 179.0619 183.9129
```

We can also use the a ‘prediction’ interval.

```
predict.lm(fit, newdata = data.frame(age=47), interval = "prediction")
```

```
##          fit      lwr      upr
## 1 181.4874 138.7833 224.1915
```

Notice the difference in the upper and lower boundaries for these predictions. The first is the prediction for the mean serum cholesterol for *individuals* of age 47 and the second is for a *single new individual* of age 47. The second prediction has to account for random variability around the mean, rather than just the precision of the estimate of the mean. This may seem like a subtle difference, but as you can see it can change our boundaries quite a bit - we need to be clear on the question we are asking.

For our multiple linear regression model explaining cholesterol as a function of age and BMI (‘mfit’), we could ask what cholesterol is predicted to be for a 60-year-old at a BMI of 21, a 60-year-old at a BMI of 26, and a 60-year-old at a BMI of 30. The standard error on the estimate of your means is obtained by setting `se.fit = TRUE`.

```
predict(mfit, newdata = data.frame(BMI = c(21,26,30), age = 60), interval = "prediction", se.fit = TRUE)
```

```
## $fit
##          fit      lwr      upr
## 1 179.2557 137.1078 221.4036
## 2 186.3885 144.3676 228.4095
## 3 192.0947 149.9339 234.2556
##
## $se.fit
```

```
##          1          2          3
## 2.025888 1.157454 2.094542
##
## $df
## [1] 397
##
## $residual.scale
## [1] 21.34293
```

Assessing the performance of the model (feedback)

Checking Residuals

Residuals are the differences between the observed response and the predicted response, and can be used to identify poorly fit data points, unequal variance (heteroscedasticity), nonlinear relationships, and examine the normality assumption.

We can plot the residuals vs x, residuals vs y, or a histogram of the residuals to see if there are any patterns. For example, plotting residuals against x (age), should be unstructured and centered at 0.

If the residuals look like they are grouped in one section of the plot, or follow a pattern, then the model is not a good fit (ie. looks quadratic - you would have a nonlinear association). If it looks like a sideways tornado, then errors are increasing with x, and this is non-constant variance.

The residuals are found in our `lm` object, 'fit', which also contains the inputs of our model; it is a list of 12.

```
str(fit)

## List of 12
## $ coefficients : Named num [1:2] 166.9 0.31
## ..- attr(*, "names")= chr [1:2] "(Intercept)" "age"
## $ residuals    : Named num [1:400] 25.13 21.27 18.24 4.55 -8.04 ...
## ..- attr(*, "names")= chr [1:400] "1" "2" "3" "4" ...
## $ effects      : Named num [1:400] -3678.3 89.45 17.61 1.86 -9.49 ...
## ..- attr(*, "names")= chr [1:400] "(Intercept)" "age" "" "" ...
## $ rank         : int 2
## $ fitted.values: Named num [1:400] 190 183 187 177 183 ...
## ..- attr(*, "names")= chr [1:400] "1" "2" "3" "4" ...
## $ assign       : int [1:2] 0 1
## $ qr           :List of 5
## ..$ qr        : num [1:400, 1:2] -20 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 ...
## .. ..- attr(*, "dimnames")=List of 2
## .. ..$ : chr [1:400] "1" "2" "3" "4" ...
## .. ..$ : chr [1:2] "(Intercept)" "age"
## .. ..- attr(*, "assign")= int [1:2] 0 1
## ..$ qraux: num [1:2] 1.05 1.02
## ..$ pivot: int [1:2] 1 2
## ..$ tol   : num 1e-07
## ..$ rank  : int 2
## ..- attr(*, "class")= chr "qr"
## $ df.residual : int 398
## $ xlevels     : Named list()
## $ call       : language lm(formula = chol ~ age, data = cholesterol)
## $ terms      :Classes 'terms', 'formula' language chol ~ age
## .. ..- attr(*, "variables")= language list(chol, age)
```

```
## ..- attr(*, "factors")= int [1:2, 1] 0 1
## ..- attr(*, "dimnames")=List of 2
## ..$ : chr [1:2] "chol" "age"
## ..$ : chr "age"
## ..- attr(*, "term.labels")= chr "age"
## ..- attr(*, "order")= int 1
## ..- attr(*, "intercept")= int 1
## ..- attr(*, "response")= int 1
## ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## ..- attr(*, "predvars")= language list(chol, age)
## ..- attr(*, "dataClasses")= Named chr [1:2] "numeric" "numeric"
## ..- attr(*, "names")= chr [1:2] "chol" "age"
## $ model      :'data.frame':  400 obs. of  2 variables:
## ..$ chol: int [1:400] 215 204 205 182 175 176 159 169 175 189 ...
## ..$ age : int [1:400] 74 51 64 34 52 39 79 38 52 58 ...
## ..- attr(*, "terms")=Classes 'terms', 'formula' language chol ~ age
## ..- attr(*, "variables")= language list(chol, age)
## ..- attr(*, "factors")= int [1:2, 1] 0 1
## ..- attr(*, "dimnames")=List of 2
## ..$ : chr [1:2] "chol" "age"
## ..$ : chr "age"
## ..- attr(*, "term.labels")= chr "age"
## ..- attr(*, "order")= int 1
## ..- attr(*, "intercept")= int 1
## ..- attr(*, "response")= int 1
## ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## ..- attr(*, "predvars")= language list(chol, age)
## ..- attr(*, "dataClasses")= Named chr [1:2] "numeric" "numeric"
## ..- attr(*, "names")= chr [1:2] "chol" "age"
## - attr(*, "class")= chr "lm"
```

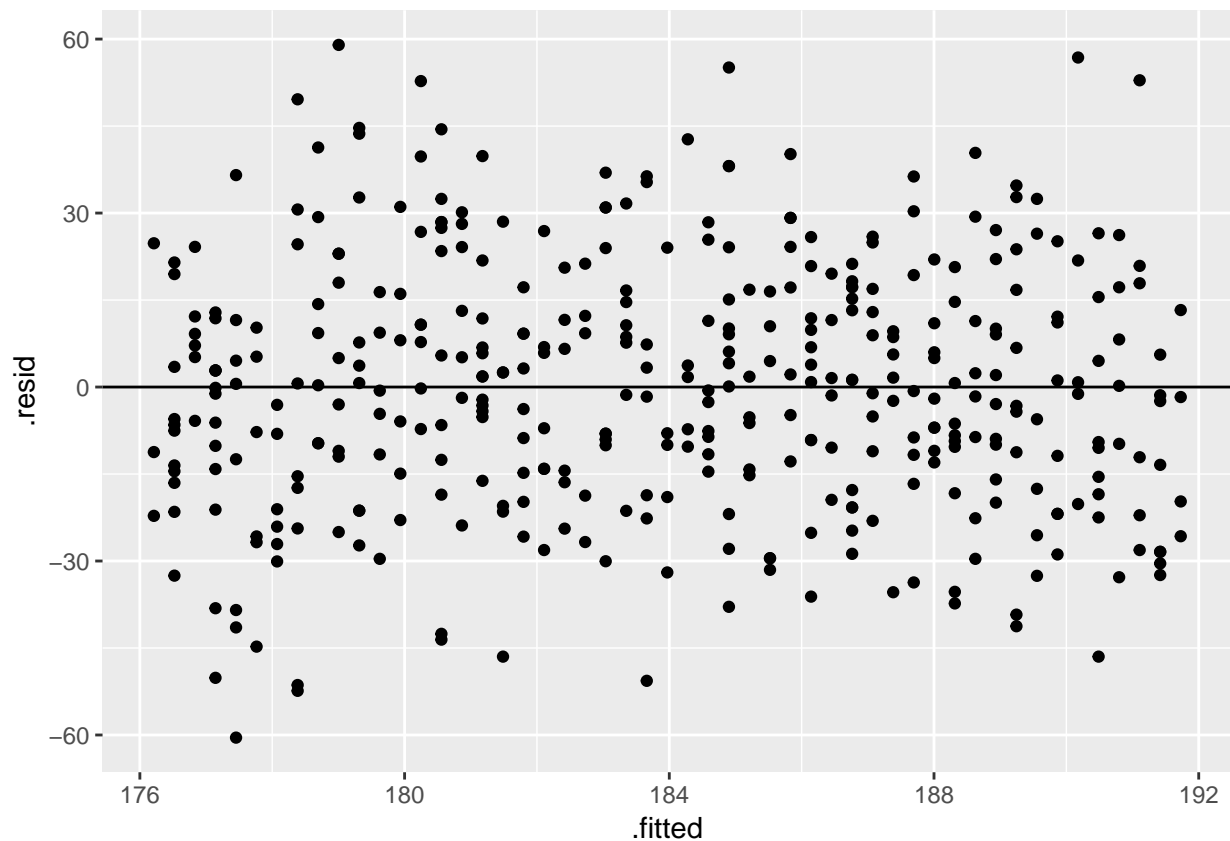
We can use the `broom()` package to get information out of linear model objects into the glorious dataframe format that we know and love. This is done using the `augment` function.

```
datfit <- augment(fit)
str(datfit)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':  400 obs. of  9 variables:
## $ chol      : int  215 204 205 182 175 176 159 169 175 189 ...
## $ age       : int   74  51  64  34  52  39  79  38  52  58 ...
## $ .fitted   : num  190 183 187 177 183 ...
## $ .se.fit   : num   1.8 1.12 1.29 1.91 1.1 ...
## $ .resid    : num  25.13 21.27 18.24 4.55 -8.04 ...
## $ .hat      : num   0.00693 0.00268 0.00351 0.00772 0.0026 ...
## $ .sigma    : num   21.7 21.7 21.7 21.7 21.7 ...
## $ .cooksd   : num   0.004717 0.001294 0.001251 0.000172 0.000179 ...
## $ .std.resid: num   1.163 0.982 0.842 0.21 -0.371 ...
```

Now that we have a data frame we can plot our residuals (`.resid`) versus our fitted data (`.fitted`).

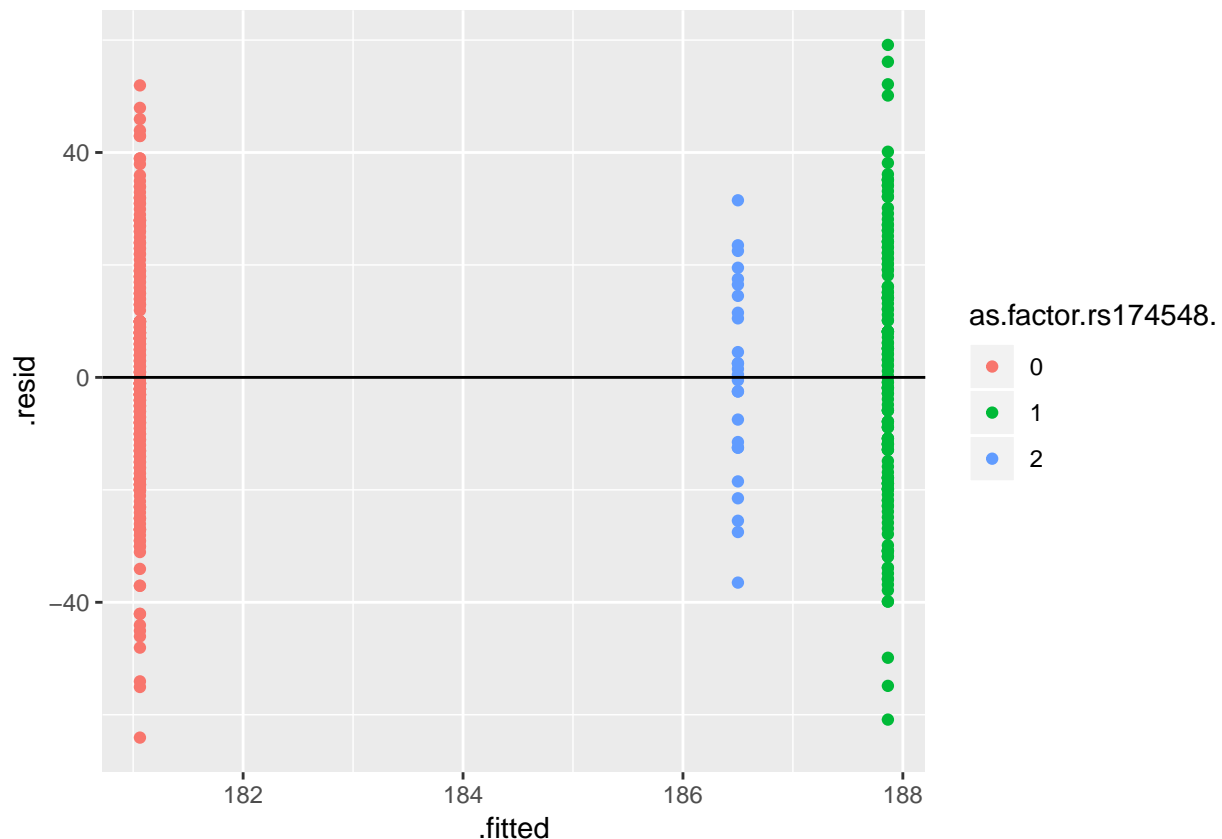
```
ggplot(datfit, aes(.fitted, .resid)) +
  geom_point() +
  geom_hline(yintercept=0, color="black")
```



You can plot the fitted and residual values with a categorical variable, but it is sometimes difficult to view patterns. For example, here is what plotting the residuals for our model of cholesterol as a function of genotype would look like.

```
anfit1 <- augment(anfit1)

ggplot(anfit1, aes(.fitted, .resid, color = as.factor(rs174548.))) +
  geom_point() +
  geom_hline(yintercept=0, color="black")
```



Instead, you can perform a statistical test of equal variance.

Bartlett's test can test whether or not population (group) variances are the same. We can see if variances are equal in our model of cholesterol as a function of sex by inputting the formula and dataset into the `bartlett.test` function.

```
bartlett.test(chol ~ factor(rs174548), data = cholesterol)
```

```
##
## Bartlett test of homogeneity of variances
##
## data: chol by factor(rs174548)
## Bartlett's K-squared = 4.8291, df = 2, p-value = 0.08941
```

The p-value is telling us that the variance is not statistically different between our populations. Our assumption of equal variance is valid.

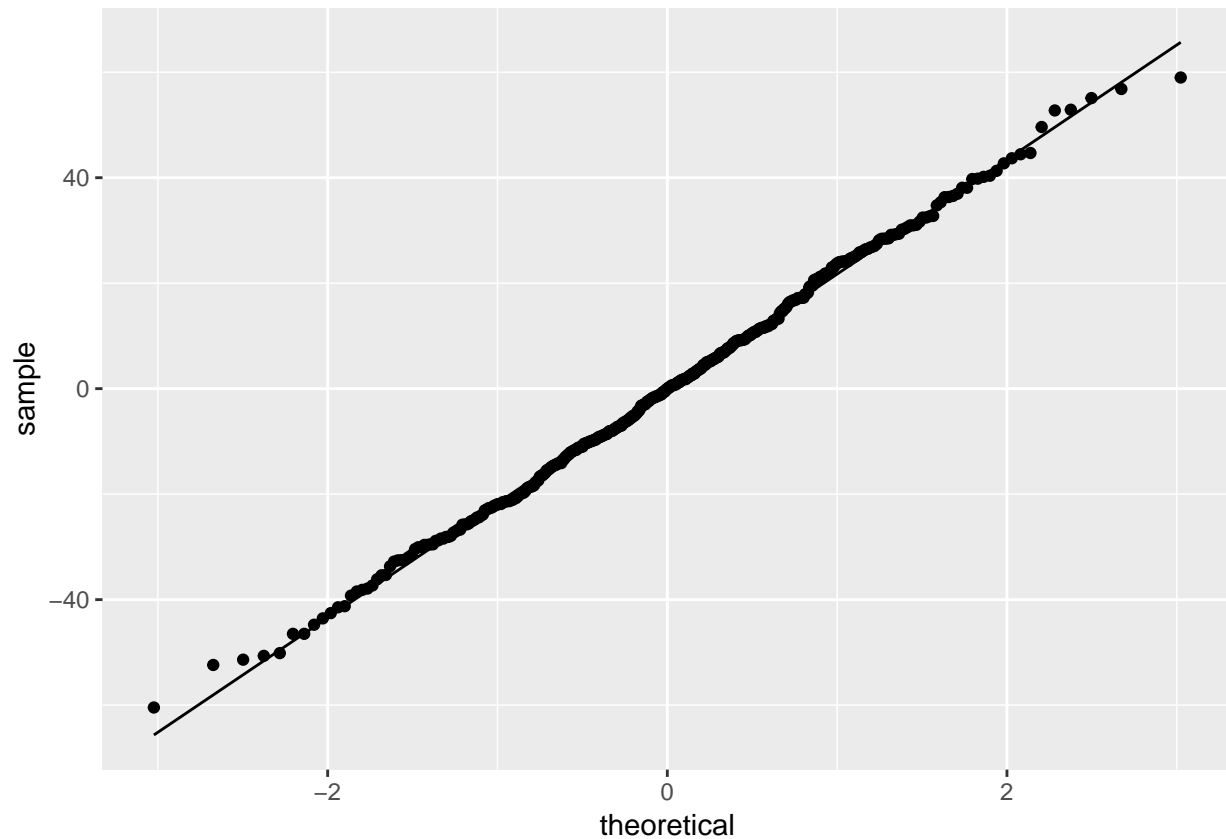
QQ-plots

QQ-plots (quantile-quantile) are a tool to answer the question: Does our data plausibly come from the (normal) distribution? The data is plotted against a theoretical distribution. Points should fall on the straight line. Any data points not fitting are moving away from the distribution.

The `stat_qq` geom from `ggplot2` allows us to plot our residuals along the y-axis in ascending order, and theoretical quantiles of a normal distribution along the x-axis. A straight line can be added to see where residuals fall with `stat_qq_line`.

```
ggplot(datfit, aes(sample = .resid)) +
  stat_qq() +
```

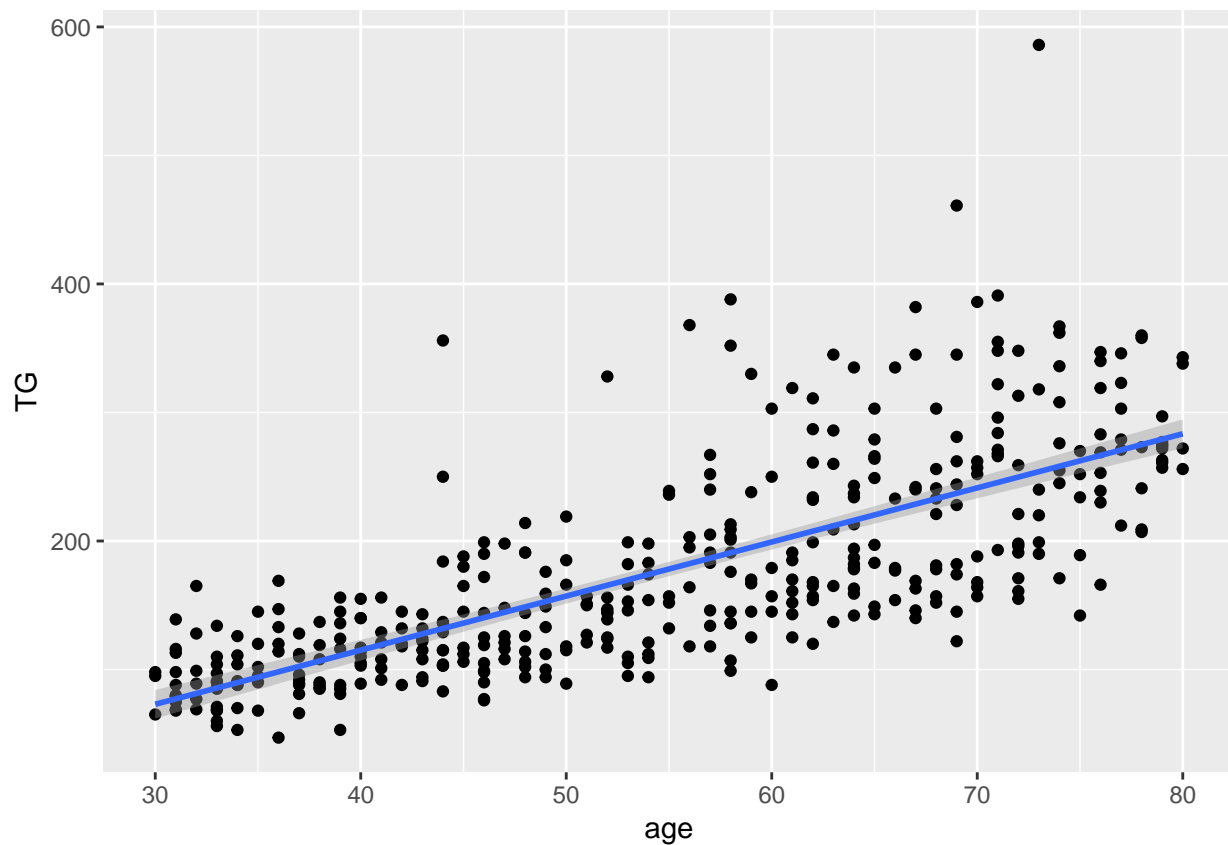
```
stat_qq_line()
```



This looks pretty straight. We likely have normality of errors.

Let's try a less perfect example and look at the relationship between age and triglycerides (TG). Make a scatterplot of age and triglycerides with a linear fit to take a look at the data.

```
ggplot(cholesterol, aes(age, TG)) +  
  geom_point() +  
  stat_smooth(method = "lm")
```



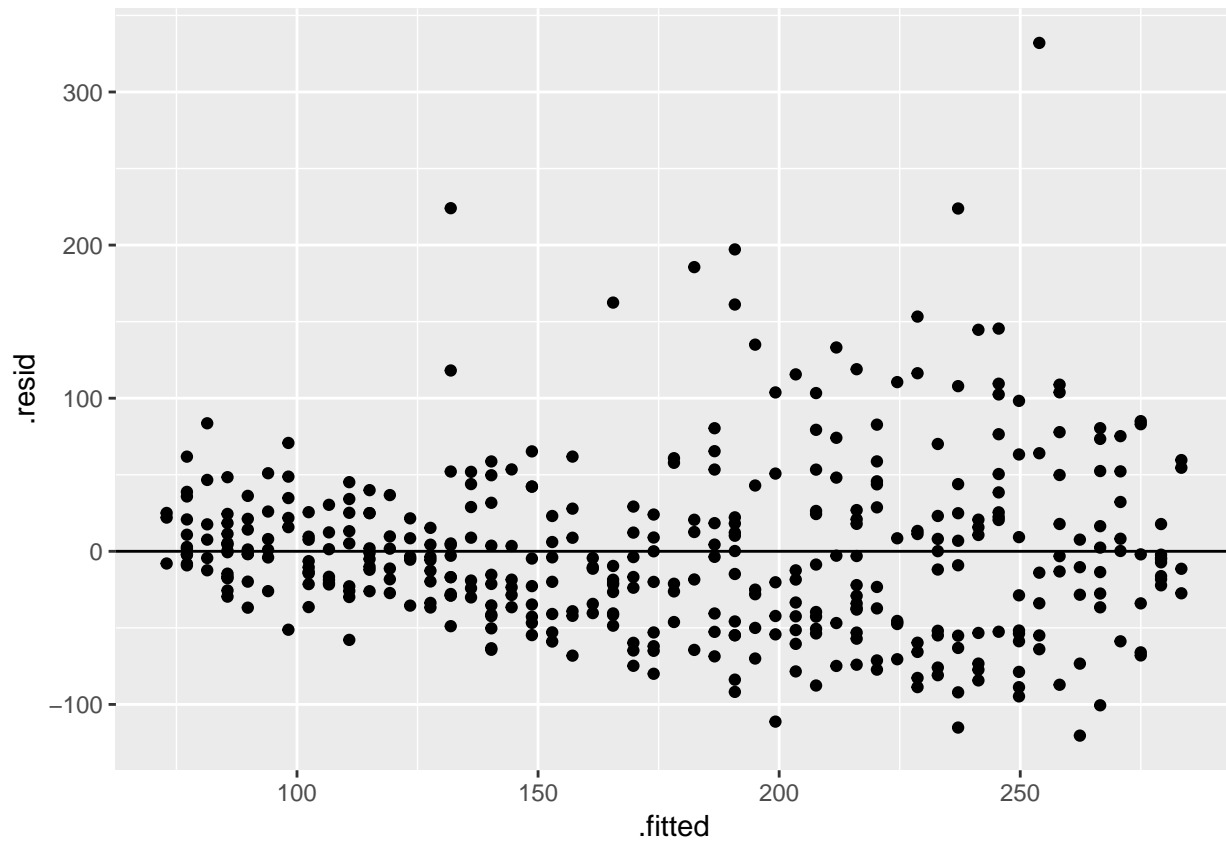
Write a linear model for triglyceride levels as a function of age. Use `broom` to get the output of the `lm` object into data frame format.

```
fitTG <- lm(TG ~ age, data = cholesterol)

datfitTG <- augment(fitTG)
```

Plot the residuals against the fitted values. Does the variance look equal across the residuals?

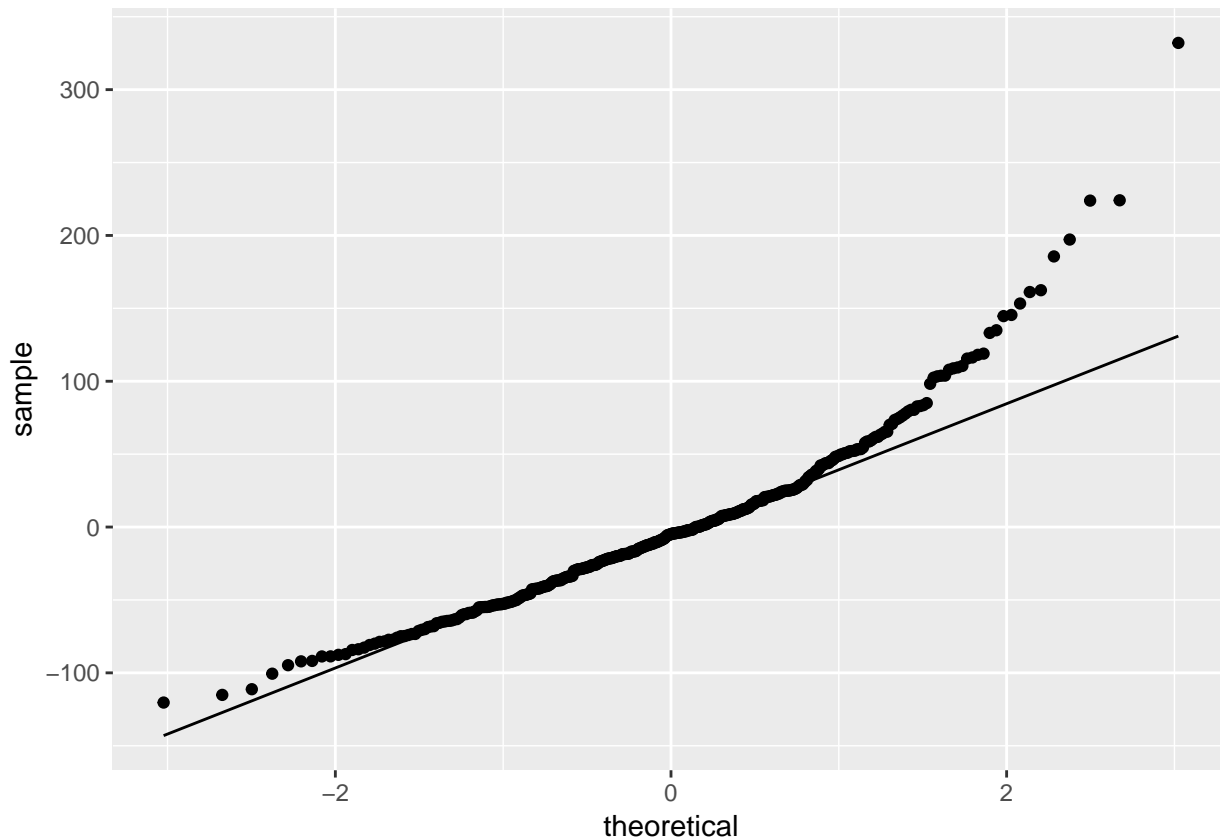
```
ggplot(datfitTG, aes(.fitted, .resid)) +
  geom_point() +
  geom_hline(yintercept=0, color="black")
```



Our residuals are increasing with increasing values of y.

What do the residuals look like in a qq-plot?

```
ggplot(datfitTG, aes(sample = .resid)) +  
  stat_qq() +  
  stat_qq_line()
```

Our qqplot points are deviating from the line suggesting a poor fit for our model.

Challenge:

For a) the anova model of the effect of the genetic factor APOE on cholesterol levels and b) the ancova model of age + APOE genotype on cholesterol levels: can you use any tools to assess whether the assumptions of your model are accurate?

Next Steps (or When Assumptions Fail)

The consequences of violating the assumptions for linear models depends, of course, on the assumption being violated. The worst offence, of course, is having non-linearity of your parameters in which case you are using the wrong model.

Our last example had a case of **non-constant variance (heteroscedasticity)**. This means that there is a mean-variance relationship (recall the tornado shape). In this case the parameter estimates are minimally impacted, however variance estimates are incorrect.

To account for this we can use:

1. Data transformation
2. Robust standard errors
3. Use a different model that does not assume constant variance (glm)

Data transformation can solve some nonlinearity, unequal variance and non-normality problems when applied to the dependent variable, the independent variable, or both. However, interpreting the results of these transformations can be tricky.

```
logfit <- lm(log(TG) ~ age, data = cholesterol)
```

```
logdat <- augment(logfit)
```

```
summary(logfit)
```

```
##
```

```
## Call:
```

```
## lm(formula = log(TG) ~ age, data = cholesterol)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
```

```
## -0.75656 -0.20390 -0.02207  0.17910  1.06931
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept) 3.7115803  0.0559237   66.37  <2e-16 ***
```

```
## age         0.0248646  0.0009866   25.20  <2e-16 ***
```

```
## ---
```

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

```
##
```

```
## Residual standard error: 0.2844 on 398 degrees of freedom
```

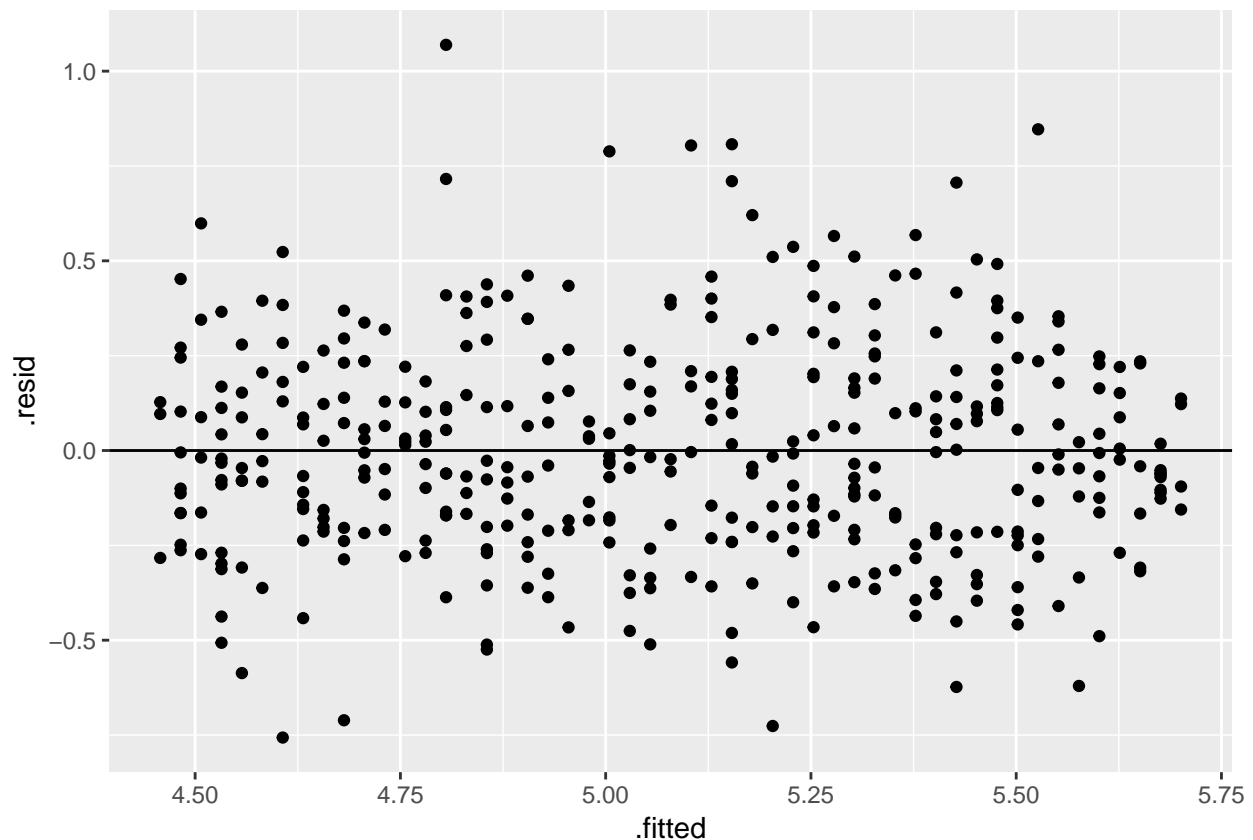
```
## Multiple R-squared:  0.6148, Adjusted R-squared:  0.6138
```

```
## F-statistic: 635.2 on 1 and 398 DF,  p-value: < 2.2e-16
```

```
ggplot(logdat, aes(.fitted, .resid)) +
```

```
  geom_point() +
```

```
  geom_hline(yintercept=0, color="black")
```



We corrected the non-constant variance issue, but it is harder to interpret our model.

Robust standard errors correctly estimate the variability of parameter estimates even under non-constant variance. This does not affect point estimates (which are minimally impacted), but corrects confidence intervals and p-values.

To do this, we use a package called **gee** (generalized estimation equation). The syntax is similar to the **lm** function, we are just adding an 'id' variable.

```
geefit <- gee(TG ~ age, data = cholesterol, id = seq(1, length(age)))
```

```
## Beginning Cgee S-function, @(#) geeformula.q 4.13 98/01/27
```

```
## running glm to get initial regression estimate
```

```
## (Intercept)          age
## -53.305930      4.208964
```

```
summary(geefit)
```

```
##
## GEE:  GENERALIZED LINEAR MODELS FOR DEPENDENT DATA
## gee S-function, version 4.13 modified 98/01/27 (1998)
##
## Model:
## Link:                               Identity
## Variance to Mean Relation: Gaussian
## Correlation Structure:              Independent
##
## Call:
## gee(formula = TG ~ age, id = seq(1, length(age)), data = cholesterol)
##
## Summary of Residuals:
##      Min       1Q   Median       3Q      Max
## -120.366372  -36.601382   -4.888487   24.529441  332.051556
##
##
## Coefficients:
##              Estimate Naive S.E.   Naive z Robust S.E.  Robust z
## (Intercept) -53.305930 11.1339178  -4.787706    8.7387366  -6.099958
## age          4.208964  0.1964165  21.428771    0.1813358  23.210880
##
## Estimated Scale Parameter:  3205.349
## Number of Iterations:  1
##
## Working Correlation
##      [,1]
## [1,]    1
```

We have the same estimates for our intercept and slope, but our error estimates have now changed. (Note: residuals in geefit are the originals.)

Use a different model that does not assume constant variance (glm)

Generalized linear models can deal with non-normal errors as well as non-linearity. They use linker functions to transform nonlinear relationships into linear ones.

Breaking the non-normality assumption will have minimal effect on estimates unless in the presence of outliers. Robust regression (ie. least squares) can be used.

Breaking the dependency assumption will have minimal effect on estimates, but the variance will be inaccurate. A different regression model for dependent data should be investigated.
