

# MATH 530/630

## Integrative Lab 2 - Linear Regression: Self Assessment

### Contents

|   |           |
|---|-----------|
| <b>Overview</b>                         | <b>2</b>  |
| <b>Logistics</b>                        | <b>2</b>  |
| <b>The Data</b>                         | <b>2</b>  |
| <b>Your Mission</b>                     | <b>3</b>  |
| <b>Exploratory Data Analysis</b>        | <b>4</b>  |
| EDA Solutions . . . . .                 | 4         |
| <b>Regression Modeling</b>              | <b>7</b>  |
| Regression Modeling Solutions . . . . . | 7         |
| <b>Residual Analysis</b>                | <b>8</b>  |
| Residual Analysis Solutions . . . . .   | 9         |
| <b>Outlier Analysis</b>                 | <b>11</b> |
| Outlier Analysis Solutions . . . . .    | 11        |
| <b>Sums of Squares</b>                  | <b>16</b> |
| Sums of Squares Solutions . . . . .     | 16        |
| <b>The NULL Model</b>                   | <b>17</b> |
| NULL Model Solutions . . . . .          | 17        |
| <b>Replicate a Plot</b>                 | <b>18</b> |
| Replicate a Plot Solutions . . . . .    | 18        |
| <b>Report your process</b>              | <b>20</b> |
| <b>Grading</b>                          | <b>20</b> |
| <b>Self-Assessment</b>                  |           |

How to:

- Start with your initial submission, and save it as a new file. Then add a sub-section to each of the sections with content called “Self-Assessment”.
- You may want to add color formatting to further highlight the self-assessment section. You can change the color of the font by using `<span style="color:deeppink"> span styles like this </span>`.
- If you want to change the color of a bulleted list, change `span` to `ul` and use the same method (pick any color).
- My code is provided here so that you can problem solve- if you need to copy and paste, do so in your self-assessment section, but you’ll need to include narrative including attribution and reflection on *what* part of that code chunk you struggled with and *why*.

## Overview

The goal of this lab is to carefully, thoroughly, and thoughtfully conduct a linear regression analysis. You are also asked to communicate clearly about the steps in your analysis process with others, by sharing your R code, output, and narrative. As such, your code cannot “stand alone”- it is meant to complement / enhance / support your narrative. This lab will be due in two stages:

1. A complete knitted `html` file submitted on Sakai.
2. A follow-up self assessment.

Using this key, your self-assessment should include even **more** narrative; where you made mistakes, you must discuss and analyze where you went wrong, and correct them without copying/pasting directly from the key (this typically means that you need to include more narrative than we provide in the key). A good self-assessment will include:

- Assessment of the accuracy and completeness of your “initial solutions”
- Correct worked solutions with some discussion and analysis of why your initial solution was incorrect, and reflection on the source of your confusion (if you got an answer correct, this is not necessary)
- Attributions as appropriate to other students who helped you, or other sources such as lecture notes, readings, online resources, etc. that helped you

## Logistics

You will use R Markdown to construct your analysis report. You’ll submit your work as an `html` file knit from your `.Rmd` file (please leave the default code chunk options for `eval = TRUE` and `echo = TRUE`). Your lab should serve as your own personal cheatsheet in the future for regression analyses. Give yourself the cheatsheet you deserve!

For all things, code and narrative, if you’re dissatisfied with a result, discuss the problem, what you’ve tried and move on (remember my 30-minute rule). You’ll need this loaded at the top:

```
library(tidyverse) # all the good stuff
library(readxl) # for reading in xlsx files
library(janitor) # for clean_names
library(knitr) # for kable
library(moderndiver) # for getting tables
library(corr) # for correlation matrix
library(skimr) # for skim
library(GGally) # for ggpairs
library(broom) # for pulling out model results
```

## The Data

You will work with an open access dataset from a publication in PLOS ONE titled: *Vitamin D Status among Thai School Children and the Association with 1,25-Dihydroxyvitamin D and Parathyroid Hormone Levels*. The data is available as an excel file on Data Dryad, where you can download the `.xlsx` file.

```
library(readxl)
path_to_xlsx <- here::here("data", "Thai_vitamin D dataset.xlsx")
vitd <- read_xlsx(path_to_xlsx, sheet = 1)
codebook <- read_xlsx(path_to_xlsx, sheet = 2) #The 2nd sheet has the codebook
```

I recommend you use the `janitor::clean_names` function, because some of these variable names start with an underscore and therefore will always need to be referenced surrounded by backticks.

```
library(janitor)
vitd <- vitd %>%
  clean_names()
glimpse(vitd)
```

Observations: 537

Variables: 19

```
$ id      <dbl> 101, 102, 104, 105, 106, 107, 108, 110, 112, 113, ...
$ sex     <chr>  "F", "F", "F", "F", "M", "M", "M", "F", "F", "F", ...
$ ageyears <dbl>  7.749487, 7.040383, 7.289528, 7.251198, 6.845995, ...
$ height  <dbl> 121.80, 120.40, 117.00, 115.60, 113.30, 121.60, 12...
$ weight  <dbl> 20.0, 18.1, 19.6, 19.3, 17.0, 23.8, 24.1, 23.1, 22...
$ bmi     <dbl> 13.48141, 12.48606, 14.31807, 14.44248, 13.24308, ...
$ zwfa    <dbl> -1.29, -1.47, -1.09, -1.17, -2.18, -0.28, -0.07, -...
$ underweight <chr> "not underweight", "not underweight", "not underwe...
$ zbfa    <dbl> -1.45, -2.24, -0.75, -0.65, -1.88, 0.28, -0.48, -1...
$ wasted  <chr> "not wasted", "wasted", "not wasted", "not wasted"...
$ zhfa    <dbl> -0.58, -0.12, -0.98, -1.20, -1.45, -0.77, 0.34, -0...
$ stunted <chr> "not stunted", "not stunted", "not stunted", "not ...
$ period  <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
$ ifyeswhen <chr> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA...
$ school  <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
$ baseline_ca <dbl> NA, NA, NA, 57.6876, NA, NA, 290.4152, 131.2554, N...
$ pth     <dbl> 66.7, 55.1, 29.5, 48.9, 20.7, 23.9, 23.4, 24.7, 33...
$ x25d    <dbl> 72.2, 73.4, 64.6, 73.6, 88.5, 85.8, NA, 75.2, 82.9...
$ x125d    <dbl> 265, 144, 320, 245, 251, 206, 200, 205, 254, 223, ...
```

That is better! You are ready to start exploring the vitd data.

## Your Mission

We'll conduct a linear regression analysis to examine the impact of age, and its interaction with sex, along with height-for-age Z score on serum 25(OH)D concentrations (a measure of vitamin D levels).

The relevant variables you will need are:

```
lm_vars <- c("sex", "_zhfa", "ageyears", "_25D")
codebook %>%
  filter(Variable %in% lm_vars) %>%
  knitr::kable()
```

Variable

Description

Unit or coding where applicable

sex

Sex of child participant

F=female, M=male

ageyears

Age of child participant

YEARS

`_zhfa`

Height-for-age z-score

Standard deviation (SD)

`_25D`

Serum 25-hydroxyvitamin D concentrations

nmol/L

Please note that the variable names in the codebook won't exactly match those in your data if you use `janitor::clean_names()`! All `_` will be replaced by `x` so that you don't have to type backticks in your later R code.

## Exploratory Data Analysis

Conduct a thorough EDA of the four variables defined above in the vitamin D dataset. Recall that a new exploratory data analysis involves three things:

- Looking at the raw values.
  - `dplyr::glimpse()`
- Computing summary statistics of the variables of interest.
  - `skimr::skim()`
  - `corrr::correlate()`
- Creating informative visualizations.
  - `ggplot2::ggplot()`
    - \* `geom_histogram()` or `geom_density()` for numeric continuous variables
    - \* `geom_bar()` or `geom_col()` for categorical variables
  - `GGally::ggpairs()`
    - \* Note that you can add transparency to points/density plots in the `aes` call, for example: `aes(colour = sex, alpha = 0.7)`

You may wish to have a level 1 header (`#`) for your EDA, then use level 2 sub-headers (`##`) to make sure you cover all three EDA bases. **At a minimum** you should answer these questions:

- How many variables/columns?
- How many rows/observations?
- Which variables are numbers?
- Which are categorical variables (numeric or character variables with variables that have a fixed and known set of possible values; aka factor variables)?
- Complete this sentence: “There is one row per...”
- What are the correlations between variables? Does each scatterplot support a linear relationship between variables? Do any of the correlations appear to be conditional on the value of a categorical variable (like `sex`)?

At this stage, you may also find you want to use `filter`, `mutate`, `arrange`, `select`, or `count`. Let your questions lead you!

## EDA Solutions

```
vitd %>%
  count(id, sort = TRUE)
```

```
# A tibble: 537 x 2
```

```
  id      n
  <dbl> <int>
1  101     1
2  102     1
3  104     1
4  105     1
5  106     1
6  107     1
7  108     1
8  110     1
9  112     1
10 113     1
```

```
# ... with 527 more rows
```

```
small_vitd <- vitd %>%
  select(ageyears, sex, x25d, zhfa)
glimpse(small_vitd)
```

```
Observations: 537
```

```
Variables: 4
```

```
$ ageyears <dbl> 7.749487, 7.040383, 7.289528, 7.251198, 6.845995, 7.7...
```

```
$ sex <chr> "F", "F", "F", "F", "M", "M", "M", "F", "F", "F", "M"...
```

```
$ x25d <dbl> 72.2, 73.4, 64.6, 73.6, 88.5, 85.8, NA, 75.2, 82.9, 7...
```

```
$ zhfa <dbl> -0.58, -0.12, -0.98, -1.20, -1.45, -0.77, 0.34, -0.33...
```

```
skim(small_vitd)
```

```
Skim summary statistics
```

```
n obs: 537
```

```
n variables: 4
```

```
-- Variable type:character -----
variable missing complete  n min max empty n_unique
sex           0         537 537  1  1    0         2
```

```
-- Variable type:numeric -----
variable missing complete  n mean  sd  p0  p25  p50  p75  p100
ageyears      0         537 537  9.89  1.69  6.16  8.38  9.95 11.26 14.03
x25d          8         529 537 74.31 15.39 36.1  64.3  73.6  84.3 127
zhfa          0         537 537 -0.98  0.84 -3.35 -1.57 -1    -0.36  1.86
```

```
hist
```

```
<U+2581><U+2587><U+2587><U+2587><U+2587><U+2587><U+2585><U+2581>
```

```
<U+2581><U+2583><U+2586><U+2587><U+2585><U+2582><U+2581><U+2581>
```

```
<U+2581><U+2582><U+2586><U+2587><U+2586><U+2583><U+2581><U+2581>
```

```
small_vitd %>%
  group_by(sex) %>%
  skim()
```

```
Skim summary statistics
```

```
n obs: 537
```

```
n variables: 4
```

group variables: sex

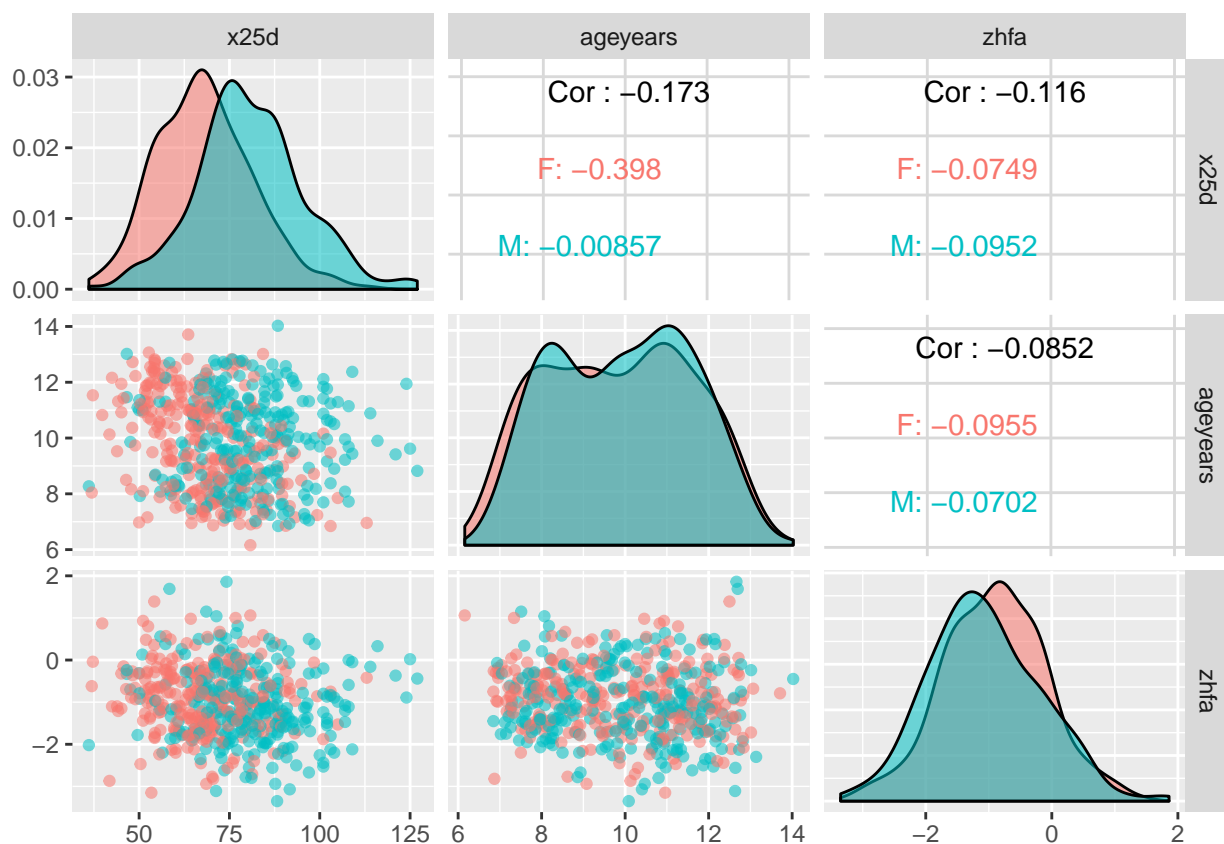
-- Variable type:numeric -----

| sex | variable | missing | complete | n   | mean  | sd    | p0    | p25   | p50   | p75   |
|-----|----------|---------|----------|-----|-------|-------|-------|-------|-------|-------|
| F   | ageyears | 0       | 271      | 271 | 9.83  | 1.74  | 6.16  | 8.3   | 9.86  | 11.18 |
| F   | x25d     | 3       | 268      | 271 | 68.08 | 13.19 | 36.9  | 58.68 | 67.35 | 76.6  |
| F   | zhfa     | 0       | 271      | 271 | -0.9  | 0.8   | -3.15 | -1.46 | -0.86 | -0.32 |
| M   | ageyears | 0       | 266      | 266 | 9.94  | 1.63  | 6.85  | 8.48  | 9.98  | 11.3  |
| M   | x25d     | 5       | 261      | 266 | 80.71 | 14.87 | 36.1  | 71.6  | 79.5  | 88.8  |
| M   | zhfa     | 0       | 266      | 266 | -1.06 | 0.88  | -3.35 | -1.67 | -1.13 | -0.47 |

p100 hist

13.71 <U+2582><U+2587><U+2587><U+2587><U+2587><U+2587><U+2586><U+2581>  
 113 <U+2581><U+2585><U+2586><U+2587><U+2585><U+2582><U+2581><U+2581>  
 1.39 <U+2581><U+2582><U+2586><U+2587><U+2587><U+2585><U+2582><U+2581>  
 14.03 <U+2583><U+2587><U+2585><U+2586><U+2587><U+2586><U+2583><U+2581>  
 127 <U+2581><U+2582><U+2583><U+2587><U+2586><U+2583><U+2581><U+2581>  
 1.86 <U+2581><U+2583><U+2586><U+2587><U+2585><U+2583><U+2581><U+2581>

```
ggpairs(small_vitd, aes(colour = sex, alpha = 0.7),
        columns = c("x25d", "ageyears", "zhfa"))
```



Note: to change color of text in a markdown list, use this command before your list: `<ul style="color:cadetblue">` and this command at the end of the list `</ul>`.

- How many variables/columns?
  - 4 (I selected just the variables I need)
- How many rows/observations?

- 537, but due to missing data, the `lm` will only use 529 observations
- Which variables are numbers?
  - `ageyears` is a number, with 0 missing values- range of 6.2 to 14
  - `x25d` is a number, with 8 missing values- range of 36 to 127
  - `zhfa` is a number, with 0 missing values- range of -3.4 to 1.9
- Which are categorical variables (numeric or character variables with variables that have a fixed and known set of possible values; aka factor variables)?
  - `sex` is a categorical variable- the codebook states that F = Female, M = Male. This means if I treat this variable as a factor, female will be the default reference level of that factor because it comes first alphabetically
- Complete this sentence: “There is one row per...”
  - There is one row per study participant
- What are the correlations between variables? Does each scatterplot support a linear relationship between variables? Do any of the correlations appear to be conditional on the value of a categorical variable (like `sex`)?
  - The correlation between vitamin d levels and age looks like it may depend on sex- there is a negative correlation for females (-0.4) but not for males (-0.0086). For girls, it seems like vitamin d levels decrease as age increases, although the linearity of this association is difficult to see from the bivariate scatterplot.
  - Height-for-age z-scores do not seem to be meaningfully correlated with vitamin d or age, for females or males.

## Regression Modeling

Fit a multiple regression model to predict serum 25(OH)D concentrations (`x25d` if you used `janitor::clean_names()`; `_25D` if you did not) and get the regression table. Your model should include:

- An intercept term
- A coefficient for age in years (`ageyears`)
- A coefficient for sex (`sex`)
- A coefficient for the interaction between age and sex (`ageyears:sex`)
- A coefficient for height-for-age z-scores (`zhfa` if you used `janitor::clean_names()`, `_zhfa` if not)

Interpret the output from the regression table (in complete sentences, but you may use bullet points to organize). You may wish to enhance the interpretability of your results by mean centering numerical predictor variables.

Some examples:

- Parallel slopes example here
- Interaction model here

The authors state: “Specifically, serum 25(OH)D concentrations were 19% higher in males at the mean age (9.9 years).” They also state: “females experienced a...4% decline in serum 25(OH)D levels for each increasing year of age; no decline was seen in male participants with increasing age”. In your narrative, walk through how you would use the numbers in the regression table to arrive at these numbers- are they accurate based on your regression model output?

## Regression Modeling Solutions

```
small_vitd <- small_vitd %>%
  mutate(age_ctr = ageyears - mean(ageyears))
```

```
dmod <- lm(x25d ~ age_ctr*sex + zhfa, data = small_vitd)
get_regression_table(dmod)
```

```
# A tibble: 5 x 7
  term      estimate std_error statistic p_value lower_ci upper_ci
  <chr>      <dbl>    <dbl>    <dbl>  <dbl>    <dbl>    <dbl>
1 intercept    66.4      1.04     63.7    0        64.3     68.4
2 age_ctr     -3.09     0.476    -6.49    0        -4.03    -2.15
3 sexM        12.5      1.18     10.6    0         10.2     14.8
4 zhfa        -1.73     0.701    -2.47   0.014     -3.10    -0.352
5 age_ctr:sexM  2.96      0.702     4.22    0         1.58     4.34
```

The modeling equation for males is:

$$\hat{y}_{male} = 78.9 - .13 * age\_ctr - 1.73 * zhfa$$

The modeling equation for females is:

$$\hat{y}_{female} = 66.4 - 3.09 * age\_ctr - 1.73 * zhfa$$

- Females are treated as the baseline for comparison for no other reason than “female” is alphabetically earlier than “male.” The  $b_{male} = 12.5$  is the vertical “bump” that males get in their vitamin d levels. Or more precisely, it is the average difference in vitamin d levels that males get relative to the baseline of females. From `skimr` output, we can see this is the difference between  $\bar{x}_F = 68$  and  $\bar{x}_M = 81$  (with some rounding going on here)
- We see that females have a lower intercept, and, as they age, they have a more steep associated average decrease in teaching scores: 3.09 vitamin d level units per year as opposed to .13 for males.
- Both males and females have the same slope for `zhfa`. In other words, in this model the associated effect of height is the same for males and females. So for every increase of one unit in height-for-age z-scores, there is on average an associated change of  $b_{zhfa} = -1.73$  (a decrease) in vitamin d levels.
- “Specifically, serum 25(OH)D concentrations were 19% higher in males at the mean age (9.9 years).”
  - Because I centered the age variable at the mean age (9.9 years), we can compare the intercepts:  $(78.9 - 66.4)/66.4 * 100 = 18.8$ ; so seems right.
- “females experienced a . . . 4% decline in serum 25(OH)D levels for each increasing year of age; no decline was seen in male participants with increasing age”
  - For each increasing year of age, females had 3.09 less units of vitamin d. If we take the mean age level of 66.4:  $3.09/66.4 * 100 = 4.65$ ; my math gives me a slightly higher number as a percentage.

## Residual Analysis

Examine the model residuals following using `get_regression_points(my_model)`. Perform a (raw) residual analysis first with a histogram, faceting by `sex`. Also look at the residuals as compared to the three predictor variables:

- $x_1$ : numerical explanatory/predictor variable of `age`
- $x_2$ : categorical explanatory/predictor variable of `sex`
- $x_3$ : numerical explanatory/predictor variable of height-for-age z-scores (`zhfa` if you used `janitor::clean_names()`, `_zhfa` if not)

Explain (a) what you are looking for in the plots, and (b) what you see in the context of assessing how “well” the linear model fits the data.

Some examples:

- Parallel slopes example here
- Interaction model here

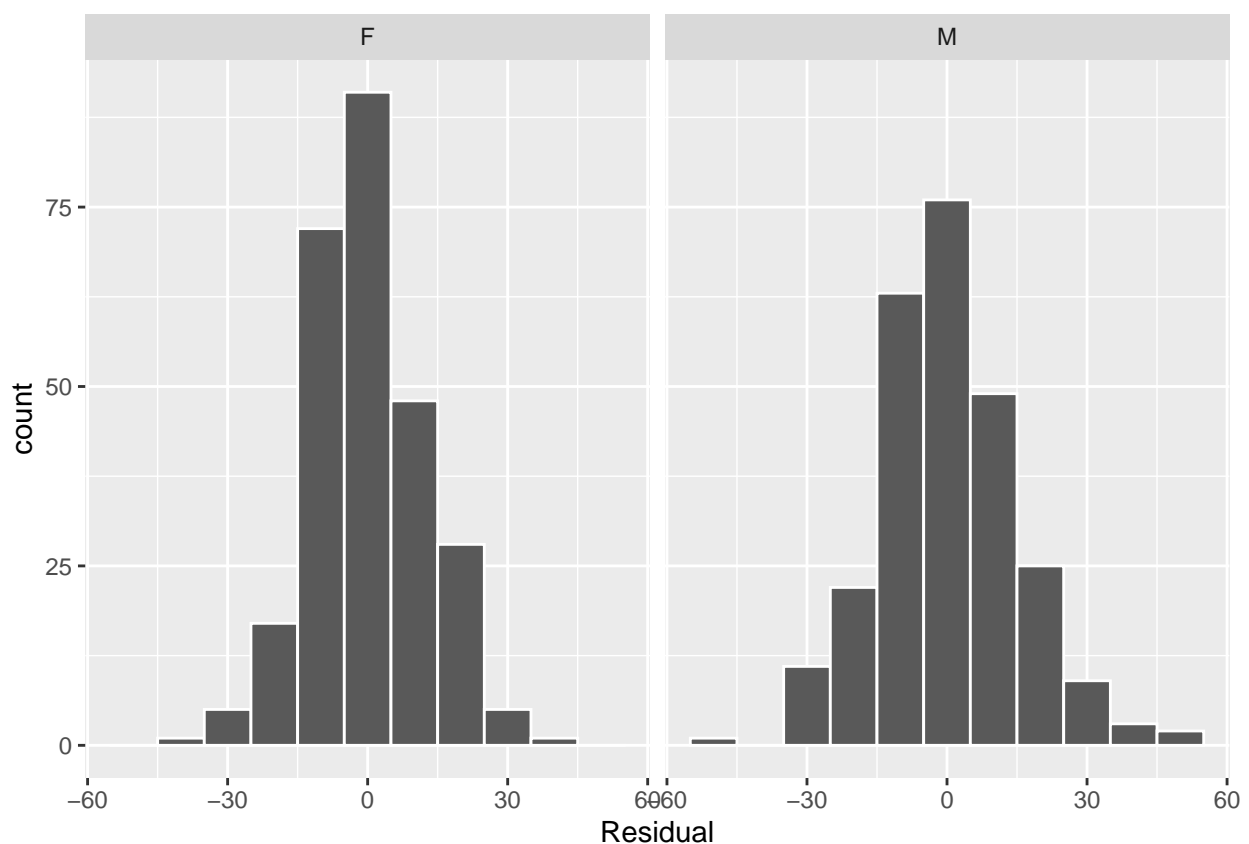


## Residual Analysis Solutions

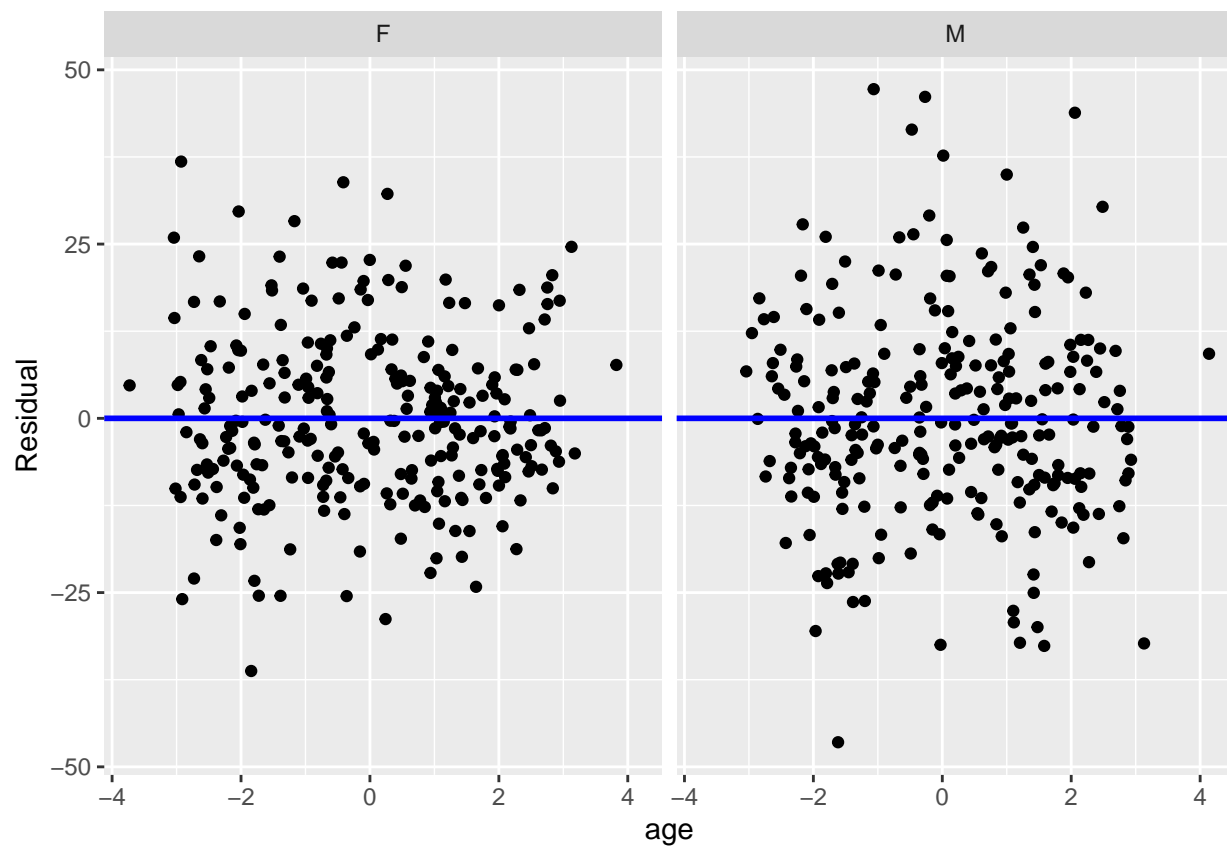
Note: to change color of text in markdown plain text, use this command before your text: `<span style="color:cadetblue">` and this command at the end `</span>`.

Residuals appear roughly unimodal, symmetric, and bell-shaped for both males and females. I don't see any distinct patterns in the residuals for either sex, although in general, larger residuals seem to be more common among males than females, suggesting that the model may work better in terms of accurately predicting vitamin d levels among females compared to males.

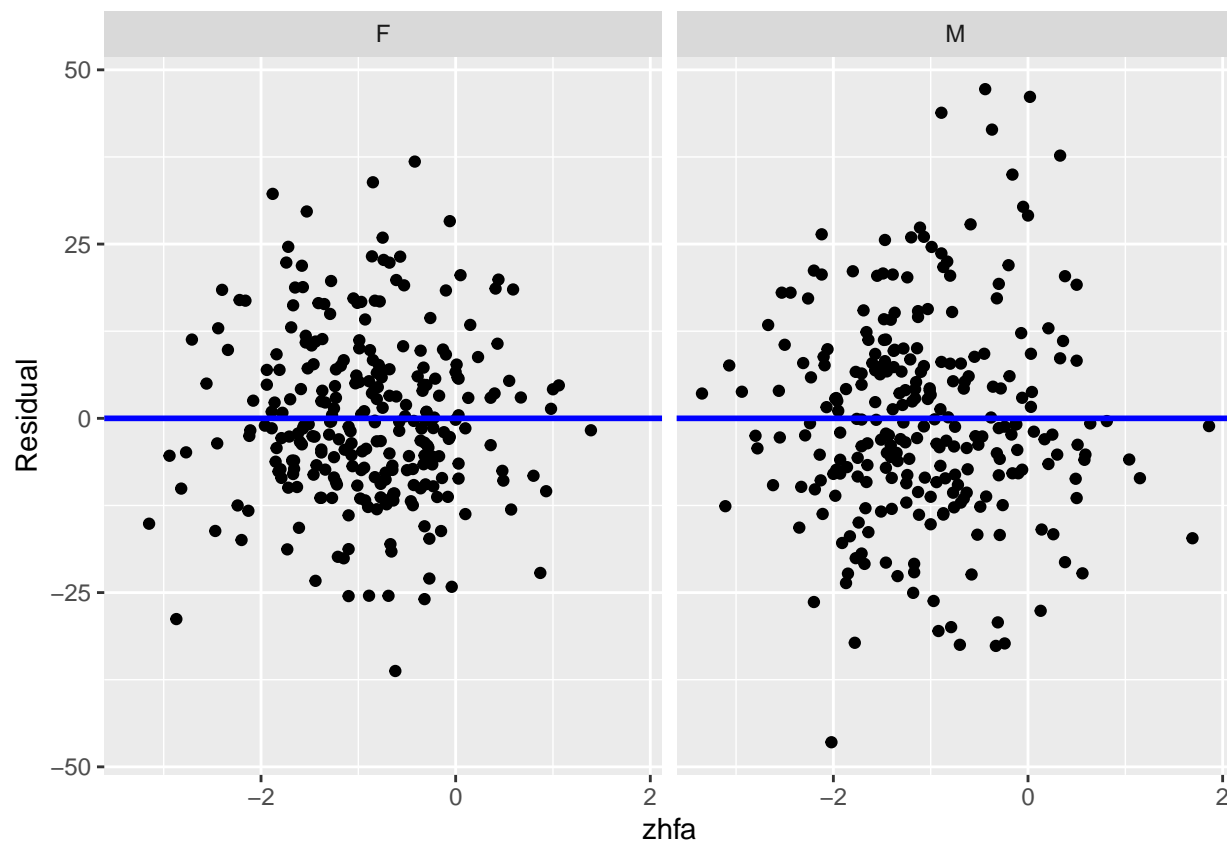
```
d_pts <- get_regression_points(dmod)
ggplot(d_pts, aes(x = residual)) +
  geom_histogram(binwidth = 10, color = "white") +
  labs(x = "Residual") +
  facet_wrap(~sex)
```



```
ggplot(d_pts, aes(x = age_ctr, y = residual)) +
  geom_point() +
  labs(x = "age", y = "Residual") +
  geom_hline(yintercept = 0, col = "blue", size = 1) +
  facet_wrap(~sex)
```



```
ggplot(d_pts, aes(x = zhfa, y = residual)) +
  geom_point() +
  labs(x = "zhfa", y = "Residual") +
  geom_hline(yintercept = 0, col = "blue", size = 1) +
  facet_wrap(~ sex)
```



## Outlier Analysis

Examine points with high leverage and discrepancy (use the externally studentized residuals as the index for discrepancy). Remember that these statistics are not in the output of the `moderndive::get_regression_points()` function- you'll want to use `broom::augment()` instead.

- How many observations would you expect to have high discrepancy in this sample?
  - You may want to account for how many females versus males you would expect
- Do you see any observations that have *both* high leverage and discrepancy?
- Do any points that are either high leverage or high discrepancy **also** have high influence on the regression estimates, as measured by Cook's distance?
- Would you exclude any observations from your model? Justify your answer either way.
- In general, would you say the overall model fit is better for females or males? Why or why not?

## Outlier Analysis Solutions

```
d_diag <- augment(dmod, data = vitd) %>%
  mutate(.ext.resid = rstudent(dmod))

# leverage
k <- 5 # Number of predictors/coefficients, i.e. intercept, age, sex, age*sex, z_score
# Or, the sum of the .hat scores is the number of coefficients, including the intercept
# k <- sum(d_diag$.hat)
```

```

mean_hat <- k/nrow(d_diag)
# Or, mean(d_diag$.hat)
d_diag %>%
  select(id, .hat) %>%
  filter(.hat > (3*mean_hat))

# A tibble: 4 x 2
  id    .hat
<dbl> <dbl>
1   303 0.0294
2   359 0.0299
3   755 0.0366
4   954 0.0390

# discrepancy
d_diag %>%
  filter(abs(.ext.resid) >= 2) %>%
  select(id, ageyears, .resid, .std.resid, .ext.resid) %>%
  arrange(desc(abs(.ext.resid)))

# A tibble: 26 x 5
  id ageyears .resid .std.resid .ext.resid
<dbl>   <dbl> <dbl>   <dbl>   <dbl>
1   616     8.82  47.2     3.51     3.55
2   156     8.27 -46.5    -3.46    -3.50
3   524     9.62  46.1     3.43     3.46
4   658    11.9  43.8     3.26     3.29
5   446     9.41  41.4     3.08     3.10
6   740     9.90  37.7     2.80     2.82
7   701     6.96  36.8     2.75     2.77
8   403     8.05 -36.2    -2.70    -2.71
9   850    10.9  35.0     2.60     2.61
10  158     9.47  33.9     2.51     2.53
# ... with 16 more rows

# leverage and discrepancy
d_diag %>%
  select(id, .hat, .ext.resid) %>%
  filter(.hat > (3*mean_hat) & abs(.ext.resid) >= 2)

# A tibble: 0 x 3
# ... with 3 variables: id <dbl>, .hat <dbl>, .ext.resid <dbl>

# discrepancy by gender
d_diag %>%
  filter(abs(.ext.resid) >= 2) %>%
  count(sex)

# A tibble: 2 x 2
  sex    n
<chr> <int>
1 F         7
2 M        19

n <- nrow(d_diag)
k <- 5 # predictors (including intercept)

```

```

d <- 4/(n - k)
d

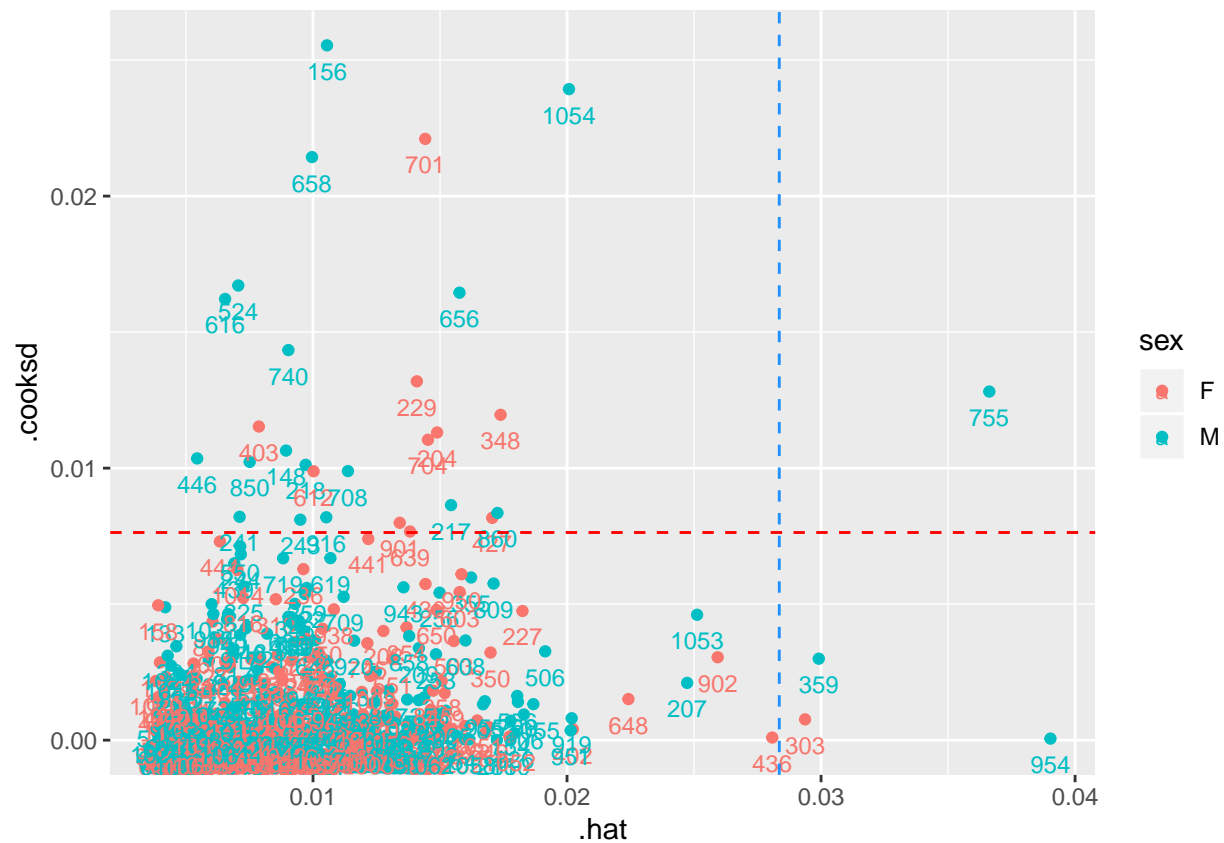
[1] 0.007633588

d_diag %>%
  filter(.cooksd > d) %>%
  select(id, x25d, .cooksd) %>%
  arrange(desc(.cooksd))

# A tibble: 28 x 3
      id  x25d .cooksd
  <dbl> <dbl>   <dbl>
1    156   36.1  0.0255
2   1054   46.6  0.0239
3    701   113  0.0221
4    658   124  0.0214
5    524   125  0.0167
6    656   109  0.0164
7    616   127  0.0162
8    740   116  0.0143
9     229   41.8  0.0132
10   755   58.4  0.0128
# ... with 18 more rows

ggplot(d_diag, aes(x = .hat, y = .cooksd, color = sex)) +
  geom_point() +
  geom_text(aes(label = id), size = 3, vjust = 2) +
  geom_hline(aes(yintercept = d), colour = "red", lty = "dashed") +
  geom_vline(aes(xintercept = 3*mean_hat), color = "dodgerblue", lty = "dashed")

```



```
ggplot(d_diag, aes(x = .ext.resid, y = .cooks, color = sex)) +
  geom_point() +
  geom_text(aes(label = id), size = 3, vjust = 2) +
  geom_hline(aes(yintercept = d), colour = "red", lty = "dashed") +
  geom_vline(xintercept = c(-2, 2), color = "dodgerblue", lty = "dashed")
```



- How many observations would you expect to have high discrepancy in this sample? You may want to account for how many females versus males you would expect.
  - DISCREPANCY: I expect 5% of observations to have high discrepancy, which is measured by externally studentized residuals, so I would expect generally  $.05 \times 529 = 26$  observations to be high on this metric. Given the roughly equal sample sizes, this means I'd expect 13 high values for males and females. Instead, I see 7 females with high leverage, compared to 19 males. This is a pretty big difference, and suggests (as I observed with the unstandardized “raw” residuals), that the model may better explain variability in vitamin d levels for females than for males.
- Do you see any observations that have *both* high leverage and discrepancy?
  - LEVERAGE: 4 observations with high leverage.
  - LEVERAGE + DISCREPANCY: 0 observations. Huzzah.
- Do any points that are either high leverage or high discrepancy **also** have high influence on the regression estimates, as measured by Cook's distance?
  - Yes! I do see one point with high leverage (`.hat`) and high cooks d values: 755 (a male)
    - \* Generally I see more green than pink above both cut-offs
  - Same story comparing high discrepancy and cooks d values: a lot of points here, mainly green. 755 doesn't show up as an extreme for discrepancy.
- Would you exclude any observations from your model? Justify your answer either way.
  - There are a lot of “flags” here, but I would not exclude anyone! If I were a vitamin d researcher, I might want to find out what the range of biologically plausible values are for vitamin d levels generally, so that I might have an objective way to identify, for example, where lab testing of vitamin d levels went awry. But as is, I have no way to identify incorrect data from unusual data.
- In general, would you say the overall model fit is better for females or males? Why or why not?
  - Across most metrics, the pattern suggests that the model explains variability in vitamin d levels better for females compared to males.

## Sums of Squares

Fill in the blanks in the following code block to calculate the Residual, Model, and Total Sums of Squares:

```
vitd_ss <- d_diag %>%
  summarise(total_ss = sum((___ - mean(___))^2),
            resid_ss = sum((___ - ___)^2),
            model_ss = sum((___ - mean(___))^2))
vitd_ss
```

Using `dplyr`, show that:

- The total sums of squares is equal to the residual plus the model sums of squares

$$total\_ss = resid\_ss + model\_ss$$

- The total sums of squares divided by  $(n - 1)$  is equal to the variance of the `y` outcome variable (*hint*: you may need to look how many observations actually contributed to the model- not the same as the original  $n$  due to missing values!)

$$var_y = \frac{total\_ss}{n - 1}$$

- The  $R^2$  value in your model output is the model sums of squares divided by the total sums of squares (*hint*: `broom::glance(my_model)`).

$$R^2 = \frac{model\_ss}{total\_ss}$$

## Sums of Squares Solutions

```
vitd_ss <- d_diag %>%
  summarise(total_ss = sum((x25d - mean(x25d))^2),
            resid_ss = sum((x25d - .fitted)^2),
            model_ss = sum((.fitted - mean(x25d))^2),
            total = resid_ss + model_ss,
            rsq = model_ss / total_ss,
            var_est = total_ss / (n()-1),
            vitd_var = var(x25d))
vitd_ss
```

```
# A tibble: 1 x 7
  total_ss resid_ss model_ss   total   rsq var_est vitd_var
  <dbl>     <dbl>   <dbl>   <dbl> <dbl>   <dbl>   <dbl>
1  125063.   95511.   29553.  125063. 0.236   237.    237.
```

```
glance(dmod)
```

```
# A tibble: 1 x 11
  r.squared adj.r.squared sigma statistic p.value    df logLik   AIC   BIC
  <dbl>      <dbl>   <dbl>   <dbl>   <dbl> <int>  <dbl> <dbl> <dbl>
1    0.236      0.230   13.5    40.5 1.33e-29     5 -2125. 4262. 4288.
# ... with 2 more variables: deviance <dbl>, df.residual <int>
```



## The NULL Model

In class, I asserted that `lm` is by default comparing the model you specify in your `lm` call to a null model defined by using a line with an intercept but slope = 0, which estimates the mean of  $y$ . Let's build an intercept-only linear regression model to prove this.

```
vitd_complete <- small_vitd %>%
  drop_na(x25d, sex, ageyears, zhfa)
int_mod <- lm(x25d ~ 1, data = vitd_complete)
get_regression_table(int_mod)

# A tibble: 1 x 7
  term      estimate std_error statistic p_value lower_ci upper_ci
<chr>      <dbl>    <dbl>    <dbl>   <dbl>   <dbl>   <dbl>
1 intercept    74.3      0.669     111.     0      73.0     75.6
```

What is the “estimate” equal to here (**hint**: look back at your `skim` output)?

```
anova(int_mod, dmod) # insert your model name from above instead of "dmod"
```

### Analysis of Variance Table

```
Model 1: x25d ~ 1
Model 2: x25d ~ age_ctr * sex + zhfa
  Res.Df    RSS Df Sum of Sq    F    Pr(>F)
1     528 125063
2     524  95511  4     29553 40.534 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Use the code above to use the `anova` function to compare the intercept-only model to your linear regression model. Look very carefully at this output and answer these questions:

- What is the RSS for line 1 (corresponding to the intercept-only model, Model 1) equal to that you calculated above?
- What is the RSS for line 2 (corresponding to Model 2) equal to that you calculated above?
- What is the Sum of Sq equal to that you calculated above?
- In < 3 sentences, explain what it means to use `lm(y ~ x + z)` versus `lm(y ~ 1)`, and what happens “under the hood” here that you now see in the `anova` output.

## NULL Model Solutions

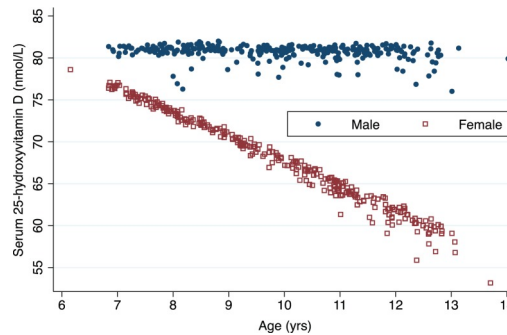
- What is the “estimate” equal to here (hint: look back at your `skim` output)?
  - The mean,  $\bar{x}_{x25d}$ , of vitamin d levels across the full sample
- What is the RSS for line 1 (corresponding to the intercept-only model, Model 1) equal to that you calculated above?
  - The total sums of squares, 125063
- What is the RSS for line 2 (corresponding to Model 2) equal to that you calculated above?
  - The residual sums of squares, 95511
- What is the Sum of Sq equal to that you calculated above?
  - The model sums of squares, 29553
- In < 3 sentences, explain what it means to use `lm(y ~ x + z)` versus `lm(y ~ 1)`, and what happens “under the hood” here that you now see in the `anova` output.
  - `lm` compares the model I entered (model 2) to an intercept-only model (model 1)- it estimates the fit of model 2 in minimizing the residual sums of squares, and compares that to what is left over or residual after using the mean vitamin d level across all children to predict individual vitamin d

levels for each child (model 1). The  $R^2$  for my model was 0.24, which is the ratio of the model sums of squares to the total sums of squares. That is the same as calculating:

$$1 - \frac{RSS_{Model2}}{RSS_{Model1}} = 1 - \frac{95561}{125063} = .236$$

- added interpretation not needed but provided here: if the two models were equivalent this number would be  $1 - 1 = 0$ . Model 2 always has to perform as least as well as Model 1. You can think of the ratio as 1 minus the proportion of variation “left over” by the null model (Model 1) that can be explained by model 2.

## Replicate a Plot



In the original paper, they presented this figure:

Figure 1 Relation between age in years and serum 25-hydroxyvitamin D (nmol/L) concentrations stratified by gender. Look at the published figure carefully:

- What appears to be the minimum value of serum 25D levels?
- What appears to be the max?
- Do those match your EDA?
- Recreate this figure in `ggplot2` using the observed data (ignore height-for-age z-scores at this point).
- Make the same plot where the y-axis should now be mapped to the fitted value for the outcome variable (using the full model with height-for-age z-scores in it). Discuss the differences you see between this plot and the previous. Were you able to recreate the published figure? Is it helpful in understanding the researchers' findings? Was it misleading at all?

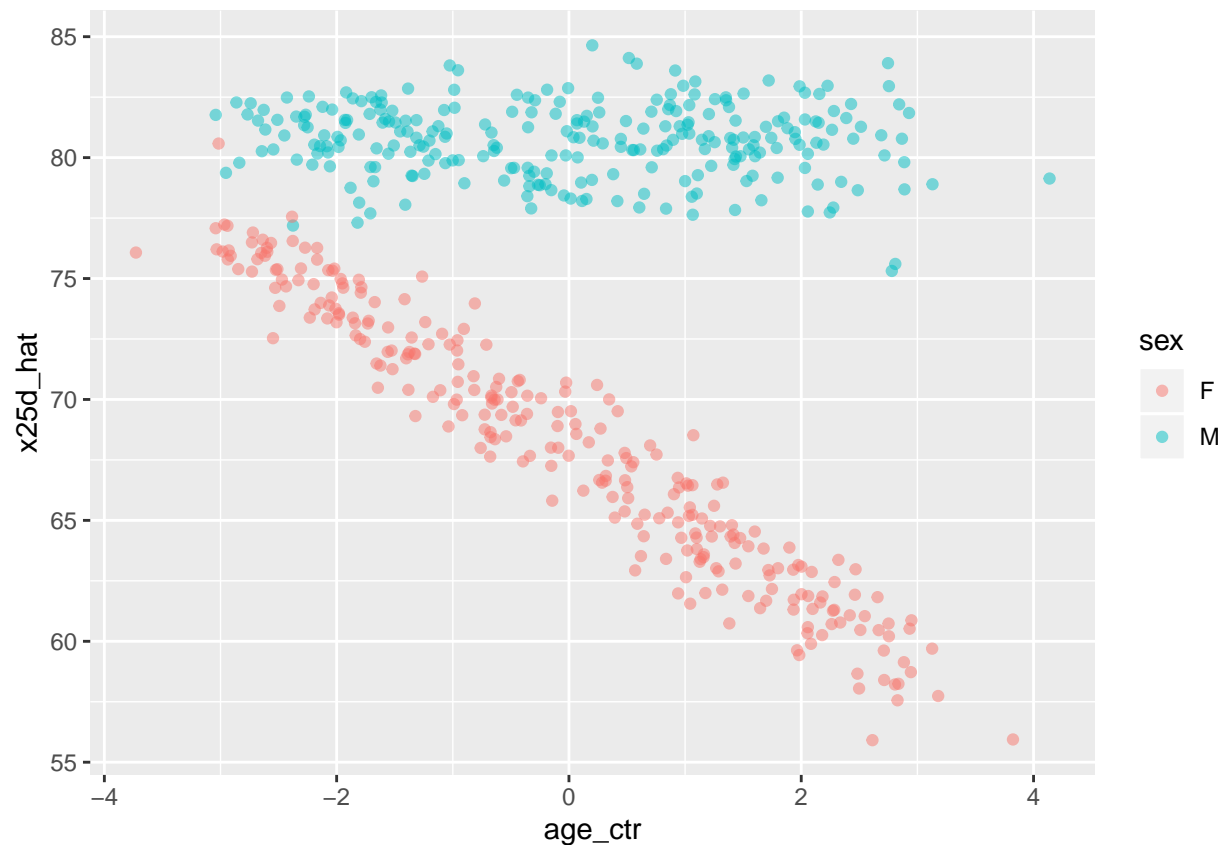
## Replicate a Plot Solutions

I'll leave the straight-forward answers as a given- they don't match the EDA, we are definitely looking at a plot of fitted values, not the original data. It may have been helpful to put that in the figure legend. The plot is helpful for visualizing the results of the regression analyses, but clarifying that the points were fitted including other variables in the model not pictured (here, `zhfa`) would have made it more reproducible and potentially less misleading. Even then though, there looks to be less variability in their fitted values (especially for males)- I'm honestly not sure how they *actually* made this plot!

```
ggplot(d_pts, aes(x = age_ctr,
                  y = x25d,
                  color = sex)) +
  geom_point(alpha = .5)
```



```
ggplot(d_pts, aes(x = age_ctr,  
                  y = x25d_hat,  
                  color = sex)) +  
  geom_point(alpha = .5)
```



## Report your process

You're encouraged to reflect on what was hard/easy, problems you solved, helpful tutorials you read, etc. Give credit to your sources, whether it's a blog post, a fellow student, an online tutorial, etc.

## Grading

This lab is worth 20 points total, scored as follows:

8 points for your initial submission being “in-good-faith”:

- 8 (Strong attempt): narrative and code reflects strong independent problem solving, with clearly thought out attempts to approach the problems and a diligent and honest effort to find solutions.
- 4 (Adequate attempt): narrative and code reflects some attempt to approach the problems, but approach appears to be superficial and lacks depth of analysis.
- 0 (No attempt or incomplete): No submission, or didn't interpret anything but left it all to the “reader”. Or more than one technical problem that is relatively easy to fix.

12 points for the quality of the final self-assessment:

- 12 (Exceptional): narrative is thorough, concise, and clearly demonstrates ability to analyze and interpret statistics as well as theoretical understanding of statistical concepts.

- 8 (Adequate): narrative addresses the questions with moderate inaccuracies in analysis and/or interpretation, or offers correct but incomplete solutions.
- 4 (Inadequate): narrative attempts to address questions with substantial inaccuracies in analysis and/or interpretation.
- 0 (Insufficient): narrative does not attempt to address questions or self-assessment is insufficient to grade