

## Ideas:

- A general (and short) introduction to `sample()`, `rnorm()`, and simulation. Reference the simulated data in C1M1 autograded.
- In the peer review assignment, we should drive home interpretations of parameters on real data, including interpretations on centered/scaled data.

In [ ]:

In [ ]:

In [ ]:

## Module 1: Peer Reviewed Assignment

### Outline:

The objectives for this assignment:

1. Learn when and how simulated data is appropriate for statistical analysis.
2. Experiment with the processes involved in simulating linear data.
3. Observe how the variance of data effects the best-fit line, even for the same underlying population.
4. Recognize the effects of standardizing predictors.
5. Interpreting the coefficients of linear models on both original and standardized data scales.

General tips:

1. Read the questions carefully to understand what is being asked.
2. This work will be reviewed by another human, so make sure that you are clear and concise in what your explanations and answers.

### A Quick Note On Peer-Reviewed Assignments

Welcome to your first peer reviewed assignment! These assignments will be a more open form than the auto-graded assignments, and will focus on interpretation and visualization rather than "do you get the right numbers?" These assignments will be graded by your fellow students (except in the specific cases where the work needs to be graded by a proctor) so please make your answers as clear and concise as possible.

```
In [3]: # This cell loads the necessary libraries for this assignment
library(tidyverse)
```

```
Warning message:
"package 'testthat' was built under R version 3.6.3"Registered S3 me
thods overwritten by 'ggplot2':
  method          from
 [.quosures       rlang
 c.quosures       rlang
 print.quosures   rlang
Registered S3 method overwritten by 'rvest':
  method          from
 read_xml.response xml2
-- Attaching packages ----- tidyve
rse 1.2.1 --
v ggplot2 3.1.1      v purrr   0.3.2
v tibble  3.0.4      v dplyr   1.0.2
v tidyr   1.1.2      v stringr 1.4.0
v readr   1.3.1      v forcats 0.4.0
Warning message:
"package 'tibble' was built under R version 3.6.3"Warning message:
"package 'tidyr' was built under R version 3.6.3"Warning message:
"package 'dplyr' was built under R version 3.6.3"-- Conflicts
----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x purrr::is_null() masks testthat::is_null()
x dplyr::lag()     masks stats::lag()
x dplyr::matches() masks tidyr::matches(), testthat::matches()
```

## Problem 1: Simulating Data

We're going to let you in on a secret. The turtle data from the autograded assignment was simulated...fake data! Gasp! Importantly, simulating data, and applying statistical models to simulated data, are very important tools in data science.

Why do we use simulated data? Real data can be messy, noisy, and we almost never *really* know the underlying process that generated real data. Working with real data is always our ultimate end goal, so we will try to use as many real datasets in this course as possible. However, applying models to simulated data can be very instructive: such applications help us understand how models work in ideal settings, how robust they are to changes in modeling assumptions, and a whole host of other contexts.

And in this problem, you are going to learn how to simulate your own data.

## 1. (a) A Simple Line

Starting out, generate 10 to 20 data points for values along the x-axis. Then generate data points along the y-axis using the equation  $y_i = \beta_0 + \beta_1 x_i$ . Make it a straight line, nothing fancy.

Plot your data (using ggplot!) with your **x** data along the x-axis and your **y** data along the y-axis.

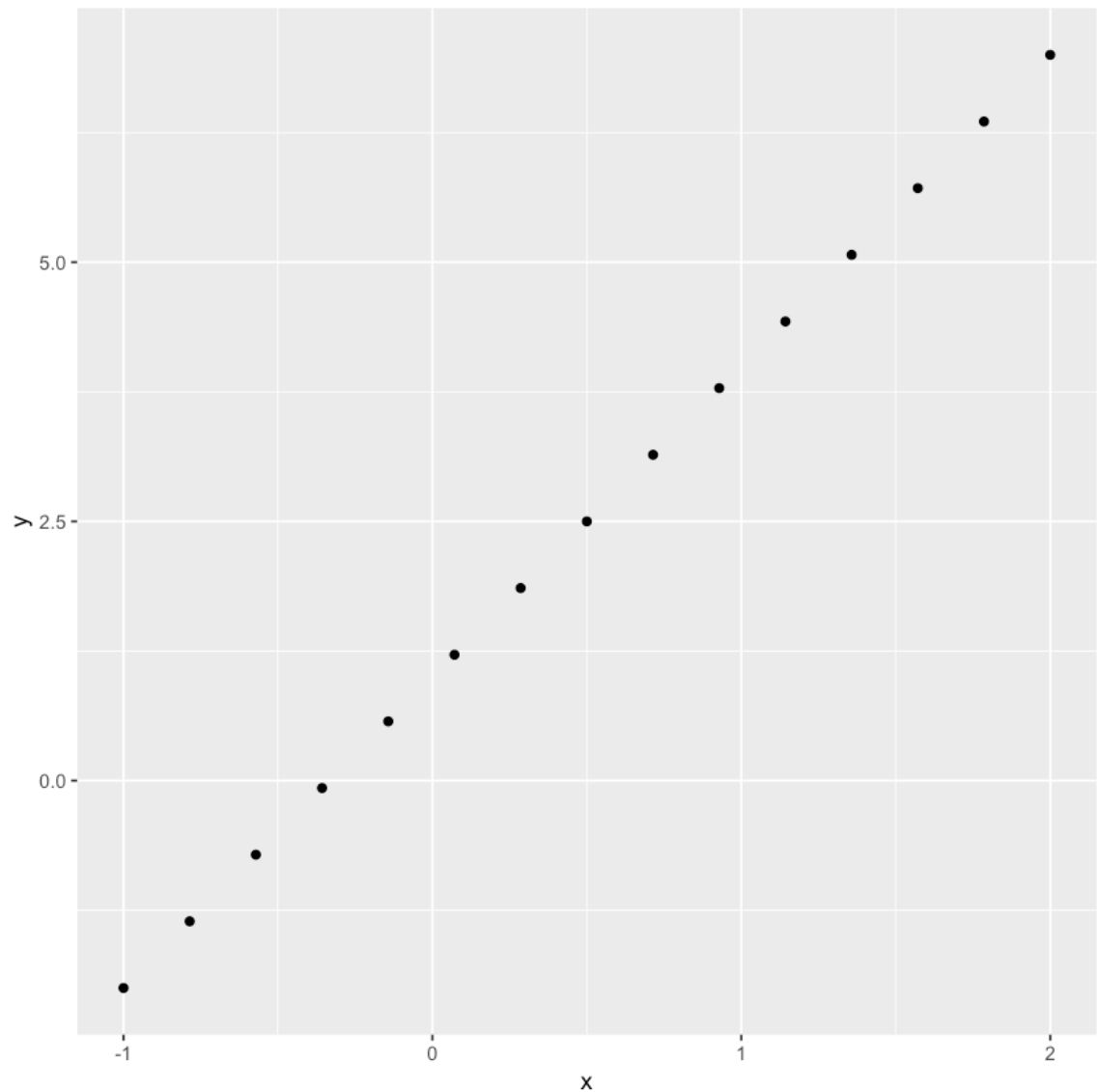
In the *Markdown* cell below the R cell, describe what you see in the plot.

**Tip:** You can generate your x-data *deterministically*, e.g., using either `a:b` syntax or the `seq()` function, or *randomly* using something like `runif()` or `rnorm()`. In practice, it won't matter all that much which one you choose.

In [30]: *# There are many possible correct solutions*

```
n = 15; x = seq(-1,2, length.out = n); b0 = 1; b1 = 3;  
y = b0 + b1*x
```

```
# ggplot  
library(ggplot2)  
df = data.frame(x = x, y = y)  
ggplot(df) +  
  geom_point(aes(x = x, y = y))
```



The data points follow the exact line  $y = 1 + 3x$ .

## 1. (b) The Error Component

That is a perfect set of data points, but that is a problem in itself. In almost any real life situation, when we measure data, there will be some error in those measurements. Recall that our simple linear model is of the form:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

Add an error term to your y-data following the formula above. Plot at least three different plots (using ggplot!) with the different values of  $\sigma^2$ .

How does the value of  $\sigma^2$  affect the final data points? Type your answer in the *Markdown* cell below the R cell.

**Tip:** To randomly sample from a normal distribution, check out the `rnorm()` function.

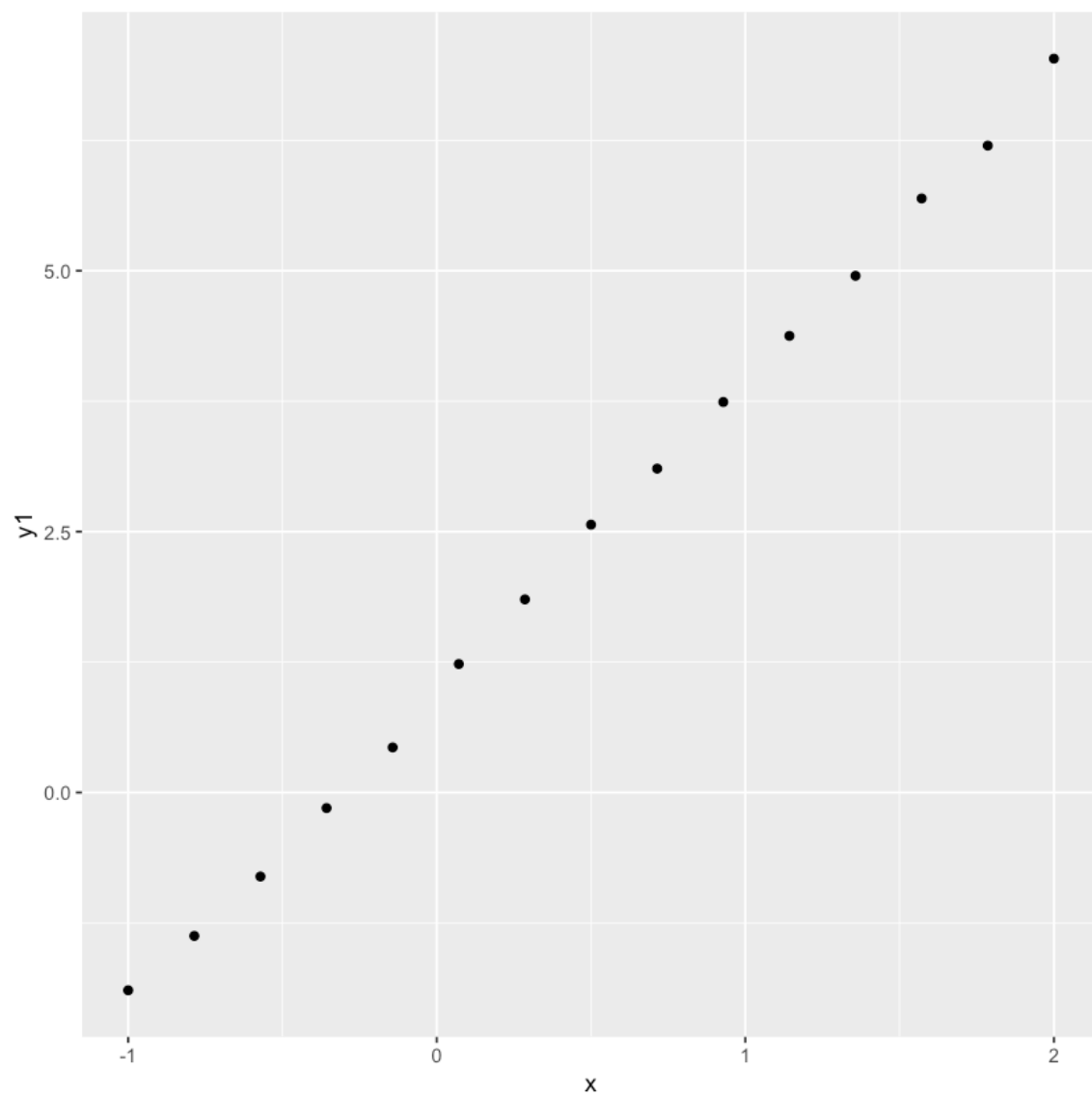
```
In [31]: # There are many possible correct solutions
set.seed(99)
e1 = rnorm(n,0, 0.1); e2 = rnorm(n,0, 1); e3 = rnorm(n,0, 10)

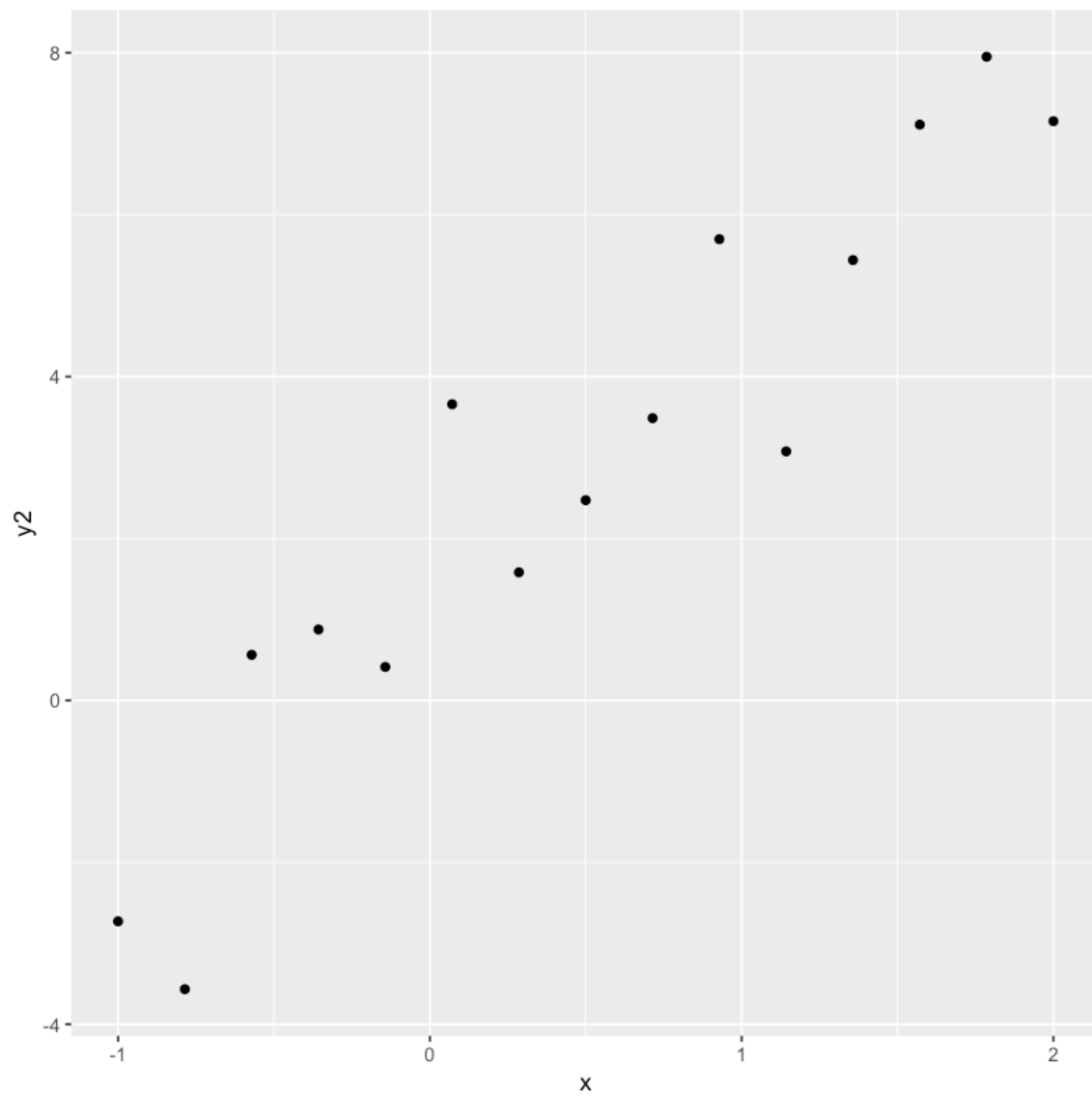
y1 = b0 + b1*x + e1
y2 = b0 + b1*x + e2
y3 = b0 + b1*x + e3

# ggplot
library(ggplot2)
df = data.frame(x = x, y1 = y1, y2 = y2, y3 = y3)
ggplot(df) +
  geom_point(aes(x = x, y = y1))

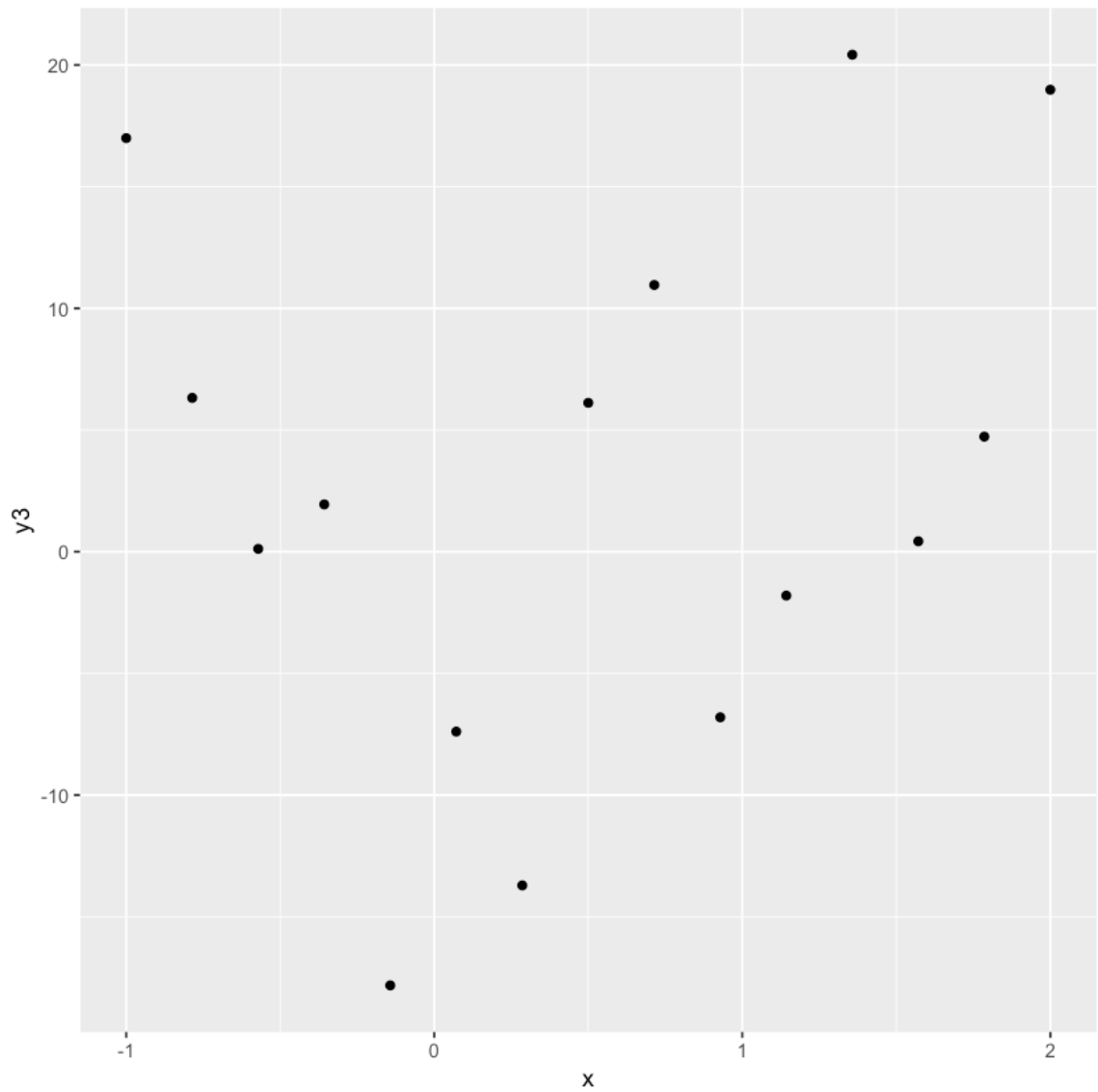
ggplot(df) +
  geom_point(aes(x = x, y = y2))

ggplot(df) +
  geom_point(aes(x = x, y = y3))
```









The larger the value of  $\sigma^2$ , the less closely the points follow a line. The last plot, with  $\sigma = 10$ , looks as if there's no linear relationship between  $x$  and  $y$ .

## Problem 2: The Effects of Variance on Linear Models

Once you've completed **Problem 1**, you should have three different "datasets" from the same underlying data function but with different variances. Let's see how those variance affect a best fit line.

Use the `lm()` function to fit a best-fit line to each of those three datasets. Add that best fit line to each of the plots and report the slopes of each of these lines.

Do the slopes of the best-fit lines change as  $\sigma^2$  changes? Type your answer in the *Markdown* cell below the R cell.

**Tip:** The `lm()` function requires the syntax `lm(y~x)` .

```
In [32]: lm1 = lm(y1 ~ x, data = df)
coef(lm1)
lm2 = lm(y2 ~ x, data = df)
coef(lm2)
lm3 = lm(y3 ~ x, data = df)
coef(lm3)
ggplot(df, aes(x = x, y = y1)) +
  geom_point() +
  geom_smooth(method = "lm")

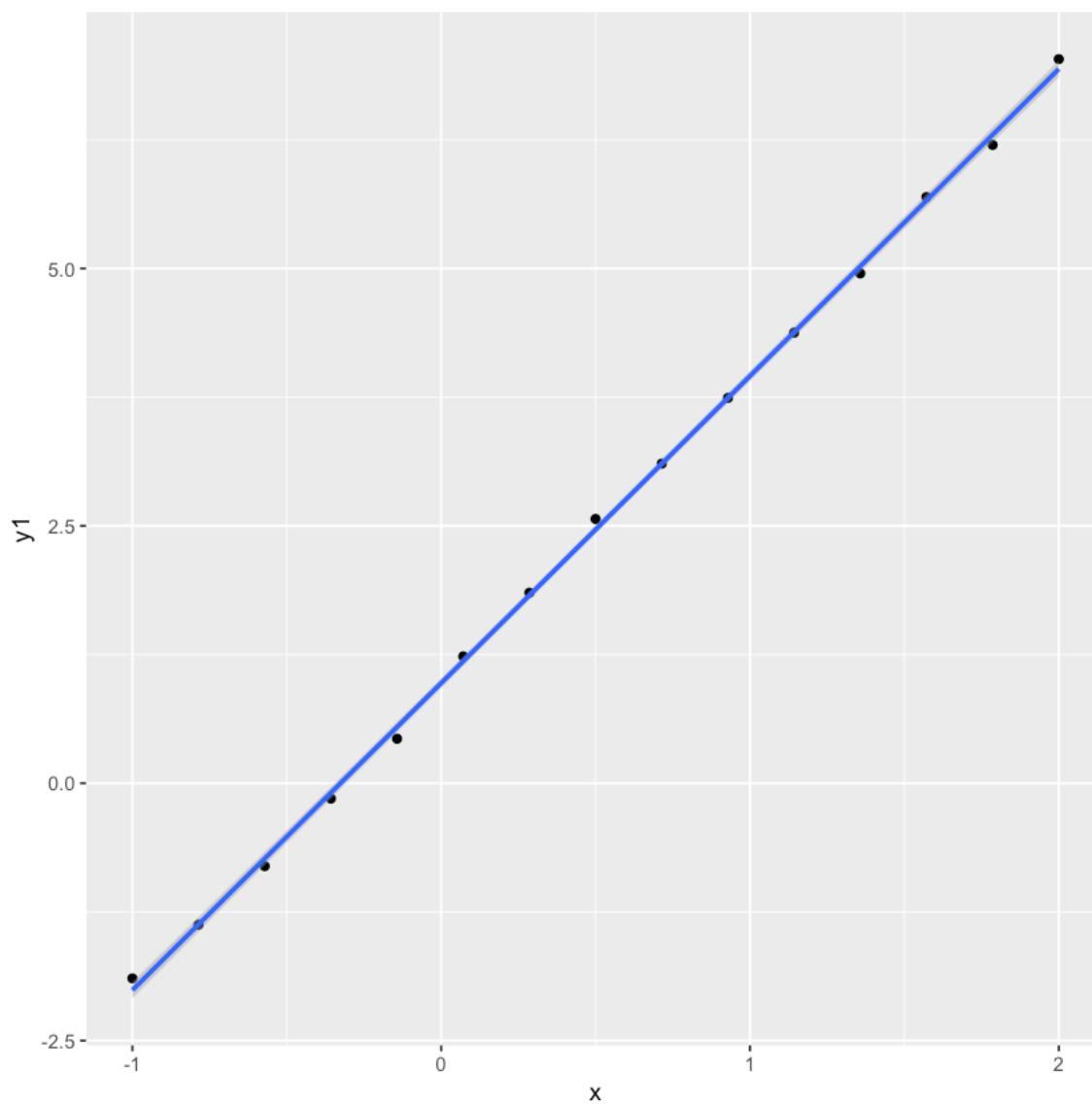
ggplot(df, aes(x = x, y = y2)) +
  geom_point()+
  geom_smooth(method = "lm")

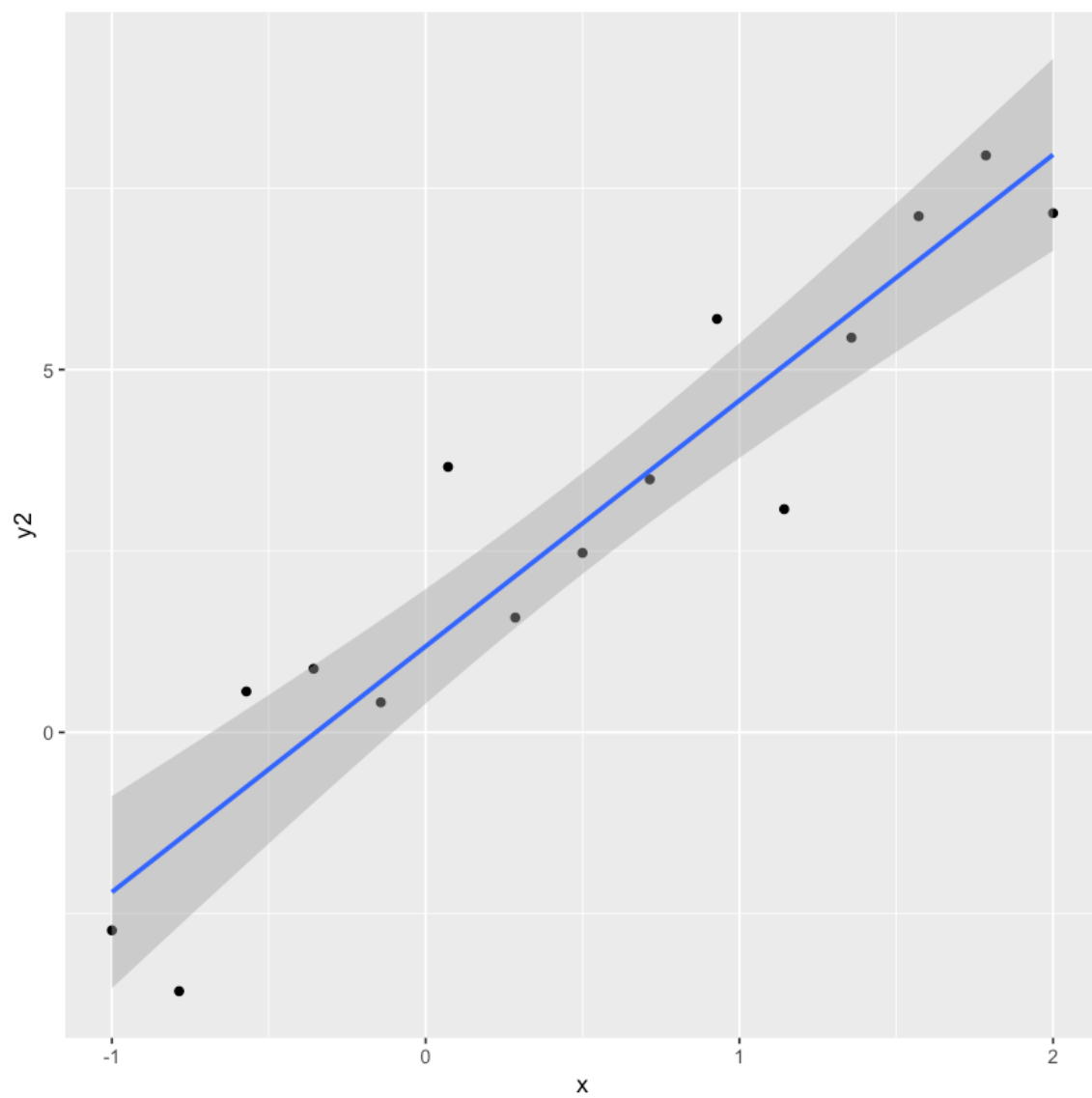
ggplot(df, aes(x = x, y = y3)) +
  geom_point()+
  geom_smooth(method = "lm")
```

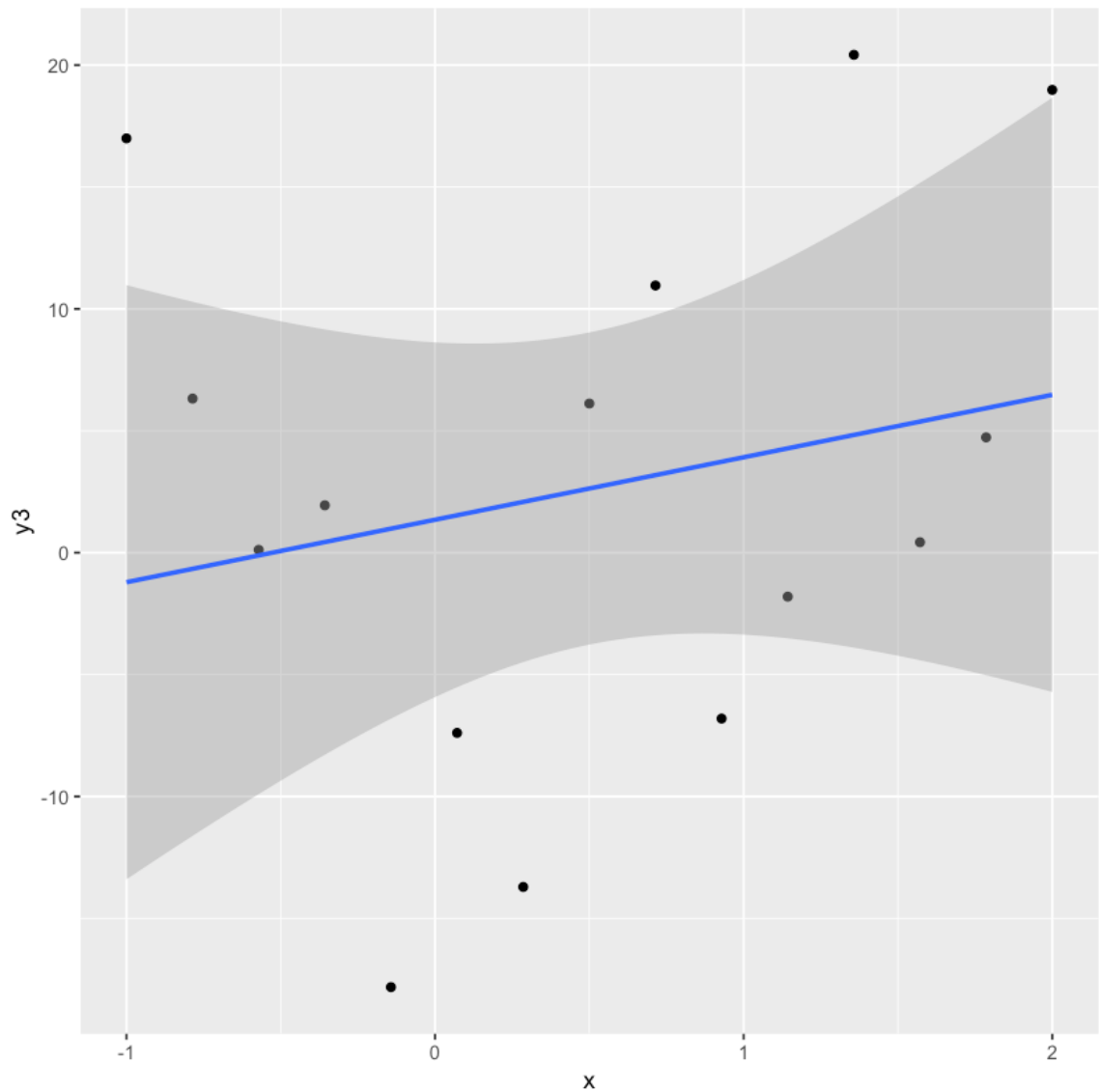
(Intercept) 0.972688967051777  
x 2.98247628818626

(Intercept) 1.18537467590712  
x 3.38796155817112

(Intercept) 1.35162740548727  
x 2.562563107789







The slopes do appear to change, especially for the highest variability case. There, for  $\sigma = 10$ , the slope is almost flat. This happens because the data points almost look like a cloud of points, i.e., no correlation.

### Problem 3: Interpreting the Linear Model

Choose one of the above three models and write out the actual equation of that model. Then in words, in the *Markdown* cell below the R cell, describe how a 1 unit increase in your predictor affects your response. Does this relationship make sense?

```
In [33]: b = coef(lm2); round(b, 2)
        yhat = b[1] + b[2]*x
```

```
      (Intercept)    1.19
              x      3.39
```

Our model is  $\hat{y} \approx 1.19 + 3.39$  (other correct/consistent rounding is OK!). On average, a one-unit increase in  $x$  will result in a 3.39 unit increase in  $y$ .

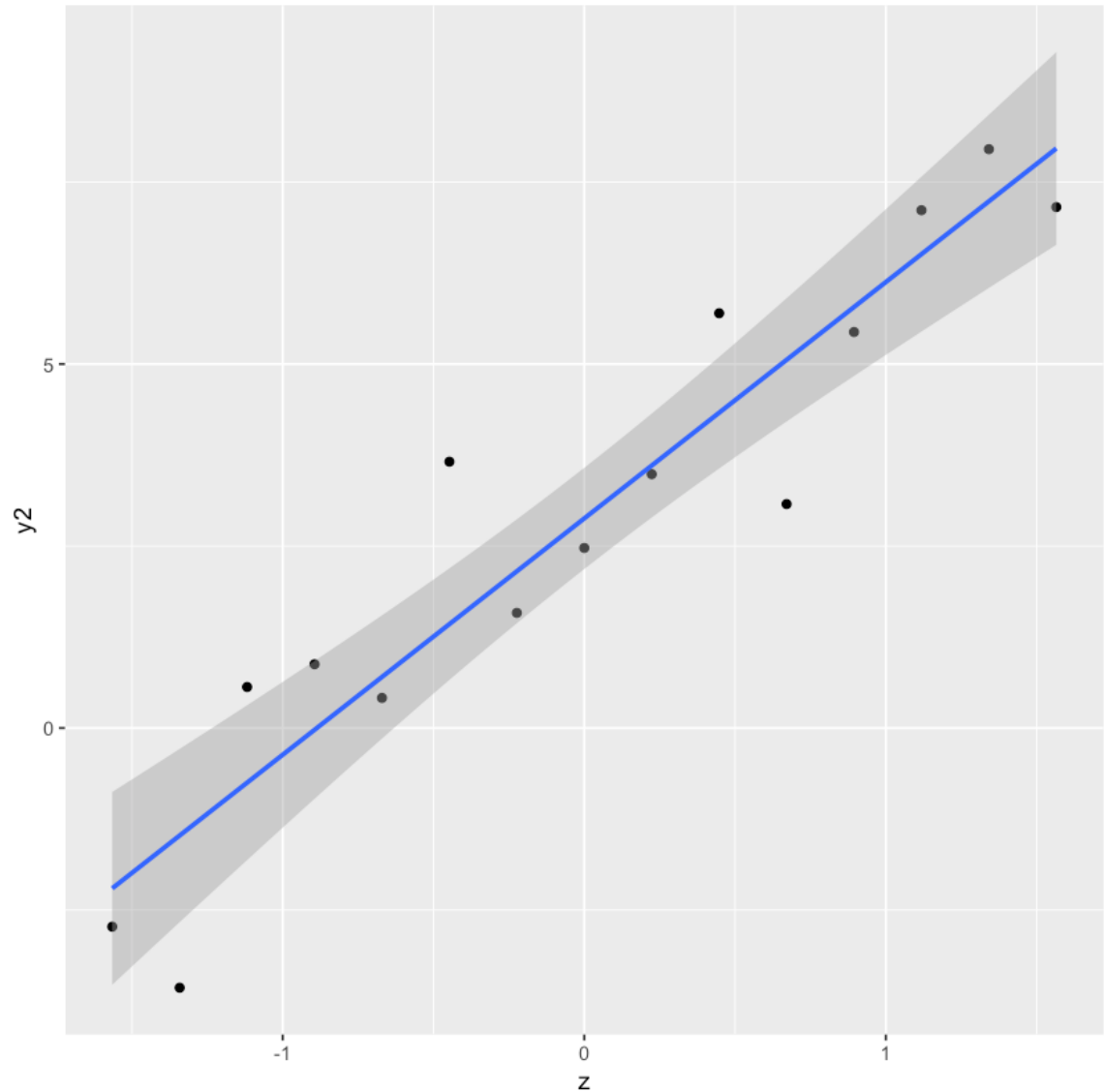
## Problem 4: The Effects of Standardizing Data

We spent some time standardizing data in the autograded assignment. Let's do that again with your simulated data.

Using the same model from **Problem 3**, standardize your simulated predictor. Then, using the `lm()` function, fit a best fit line to the standardized data. Using `ggplot`, create a scatter plot of the standardized data and add the best fit line to that figure.

```
In [34]: z = (x - mean(x))/sd(x)
df$z = z
lm_standard = lm(y2 ~ z, data = df)

ggplot(df, aes(x = z, y = y2)) +
  geom_point() +
  geom_smooth(method = "lm")
```



## Problem 5: Interpreting the Standardized Model



Write out the expression for your standardized model. In words, in the *Markdown* cell below the R cell, describe how a 1 unit increase in your standardized predictor affects the response. Is this value different from the original model? If yes, then what can you conclude about interpretation of standardized predictors vs. unstandardized predictors.

```
In [35]: b_standard = coef(lm_standard); b_standard
        yhat2 = b_standard[1] + b_standard[2]*z
```

```
      (Intercept)  2.87935545499268
              z    3.24673386395431
```

Our model is  $\hat{y} \approx 2.88 + 3.25$  (other correct/consistent rounding is OK!). On average, a one-unit increase in  $z$  - that is, a one-standard deviation increase in  $x$  - will result in a 3.25 unit increase in  $y$ .