

transform_lab_sol

2025-10-10

This non-graded lab will focus heavily on linear regression with linear-linear, log-linear, and log-log all making an appearance. The last section is dedicated to indicators.

NOTE: I want you to go through Question 1-13 first. You may then either do the transformation section next or jump directly to the indicators section (starting at Q24). Do whichever you feel like you'd like the most in-class lab help with.

The data set we will be using today is actually a super fun one an old mentor of mine collected on...LEGOs! I use to love LEGOs growing up so now you get to play with LEGOs (...data set).

First, let's read in the data and look at the first few rows

```
legos <- read.csv('https://vinnys-classes.github.io/data/legos_data.csv')
legos$Year <- as.factor(legos$Year) #RUN THIS!!
head(legos)
```

##	Item_Number	Set_Name	Theme	Pieces	Year	Pages	Minifigures
## 1	10859	My First Ladybird Duplo		6	2018	9	NA
## 2	10860	My First Race Car Duplo		6	2018	9	NA
## 3	10862	My First Celebration Duplo		41	2018	9	NA
## 4	10864	Large Playground Brick Box Duplo		71	2018	32	2
## 5	10867	Farmers' Market Duplo		26	2018	9	3
## 6	10870	Farm Animals Duplo		16	2018	8	NA

##	Packaging	Unique_Pieces	Size	amazon_price	age
## 1	Box	5	Large	16.00	1
## 2	Box	6	Large	9.45	1
## 3	Box	18	Large	39.89	1
## 4	Plastic box	49	Large	56.69	2
## 5	Box	18	Large	36.99	2
## 6	Box	13	Large	9.99	2

The variables are...

- 1) Item_Number: ID
- 2) Set_Name: The selling name of the lego set
- 3) Theme: One of three themes
- 4) Pieces: Number of pieces in the set
- 5) Year: Year the set was made as a nominal categorical variable (called factor's in R)
- 6) Pages: Number of pages in the booklet
- 7) Minifigures: Number of "people" sold with the set

- 8) Package: What type of packaging the set comes in
- 9) Unique Pieces: How many unique lego blocks are in the set
- 10) Size: The size of the blocks, with two levels
- 11) amazon_price: Price of the on Amazon as of a few years ago
- 12) age: the lowest age the company recommends for the data set

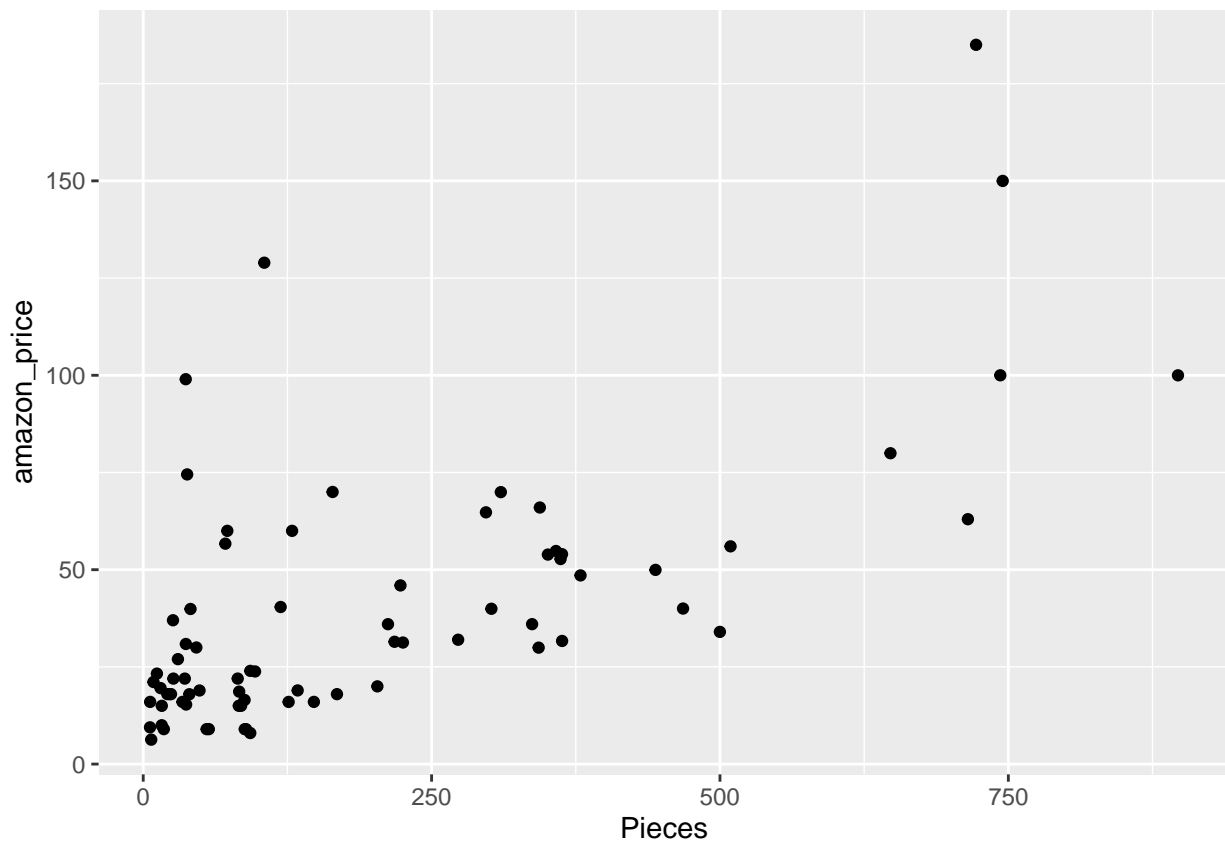
Linear Regression

Q1

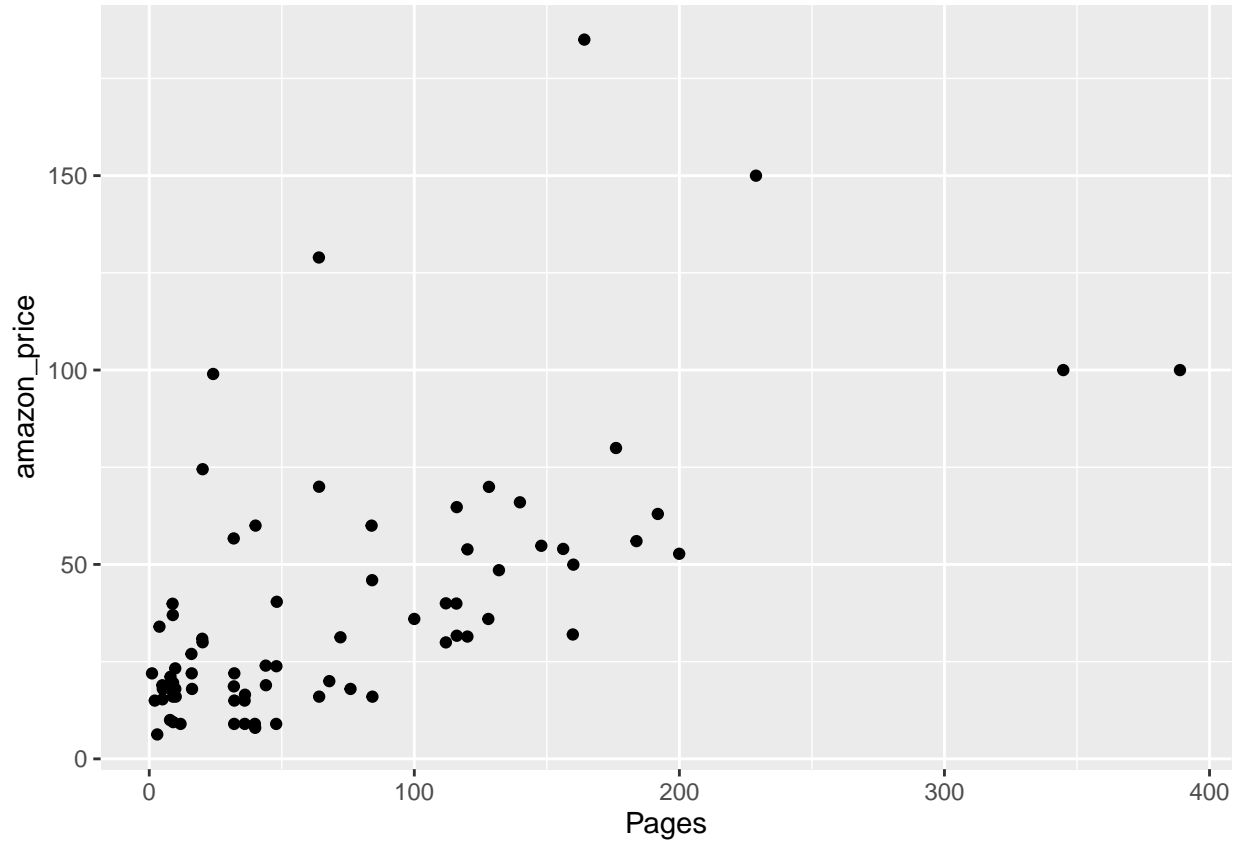
Please make three scatterplots. All three should have amazon_price as the y-axis and the three x-axis should be the variables Pieces, Pages, and Minifigures.

```
library(ggplot2)

ggplot(data = legos,
       aes(x = Pieces,
           y = amazon_price)) +
  geom_jitter(width = .2)
```

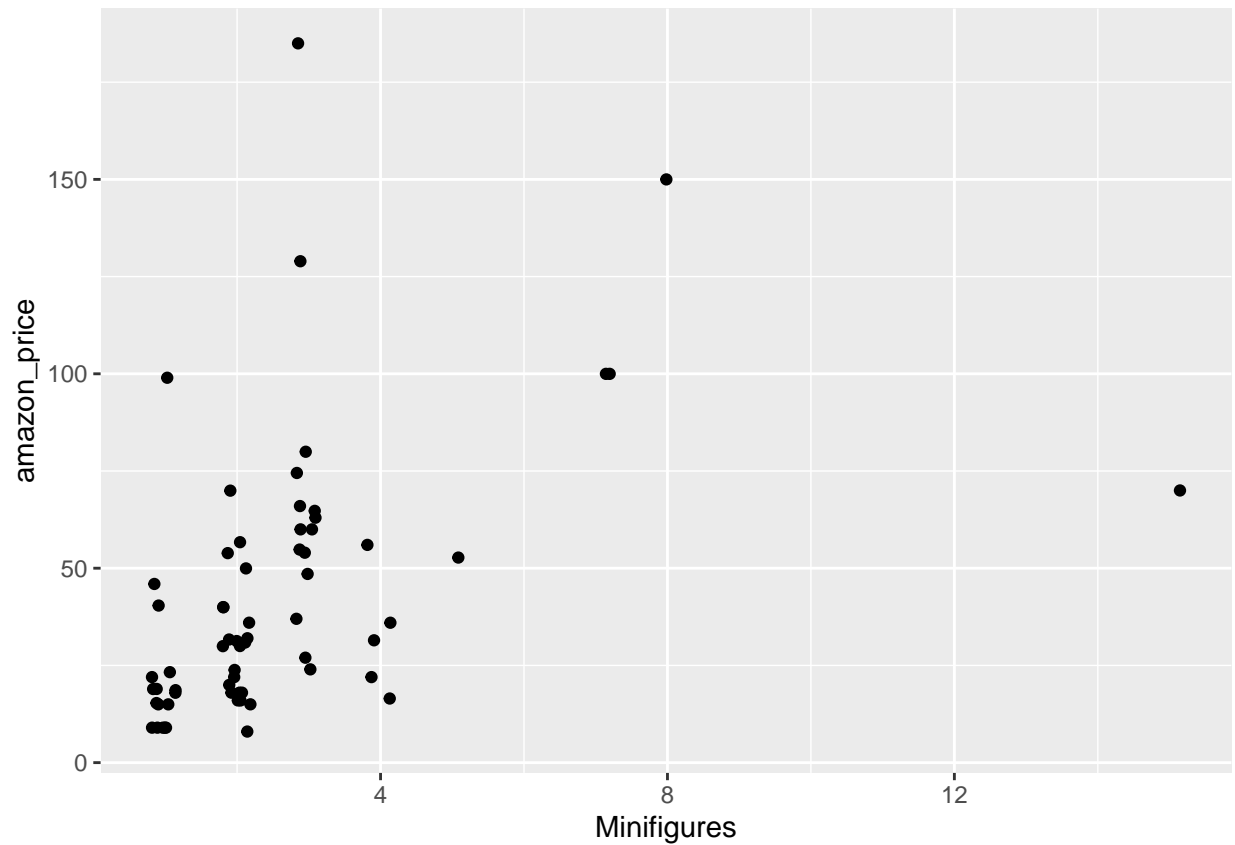


```
ggplot(data = legos,
       aes(x = Pages,
           y = amazon_price)) +
  geom_jitter(width = .2)
```



```
ggplot(data = legos,
       aes(x = Minifigures,
           y = amazon_price)) +
  geom_jitter(width = .2)
```

```
## Warning: Removed 10 rows containing missing values or values outside the scale range
## ('geom_point()').
```

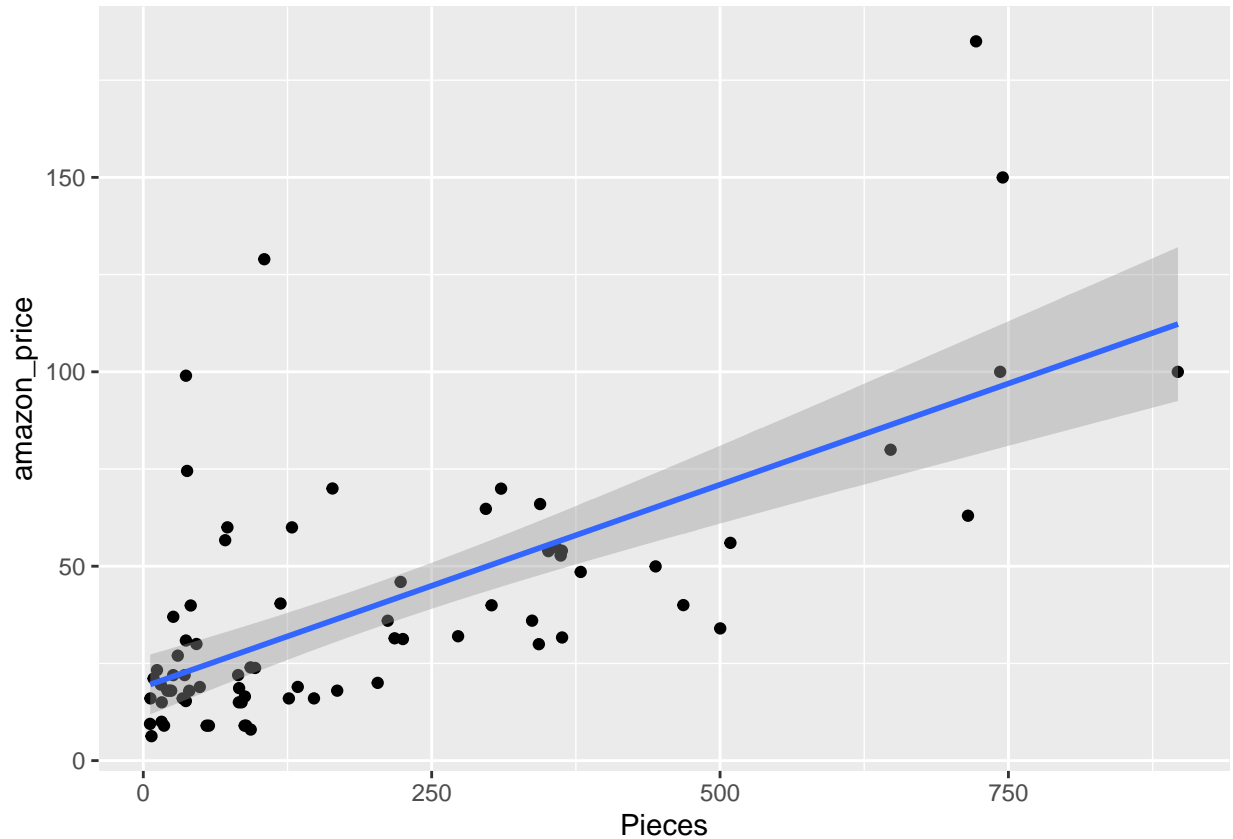


Q2

Using `geom_smooth()`, please plot a best-fit-line (by using the 'lm' method of `geom_smooth`) to the scatterplot of amazon price by number of pieces. Describe the scatterplot by noting it's direction, form, outliers, and strength please.

```
ggplot(data = legos,
       aes(x = Pieces,
           y = amazon_price)) +
  geom_jitter(width = .2) +
  geom_smooth(method = 'lm')
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



Q3

Using the `lm()` function, please fit a linear model with amazon price as the response variable and the number of pieces as the explanatory variable. Print out the summary of the model using the `summary()` function

```
my_line_mod <- lm(amazon_price ~ Pieces,
                  data = legos)
```

Q4

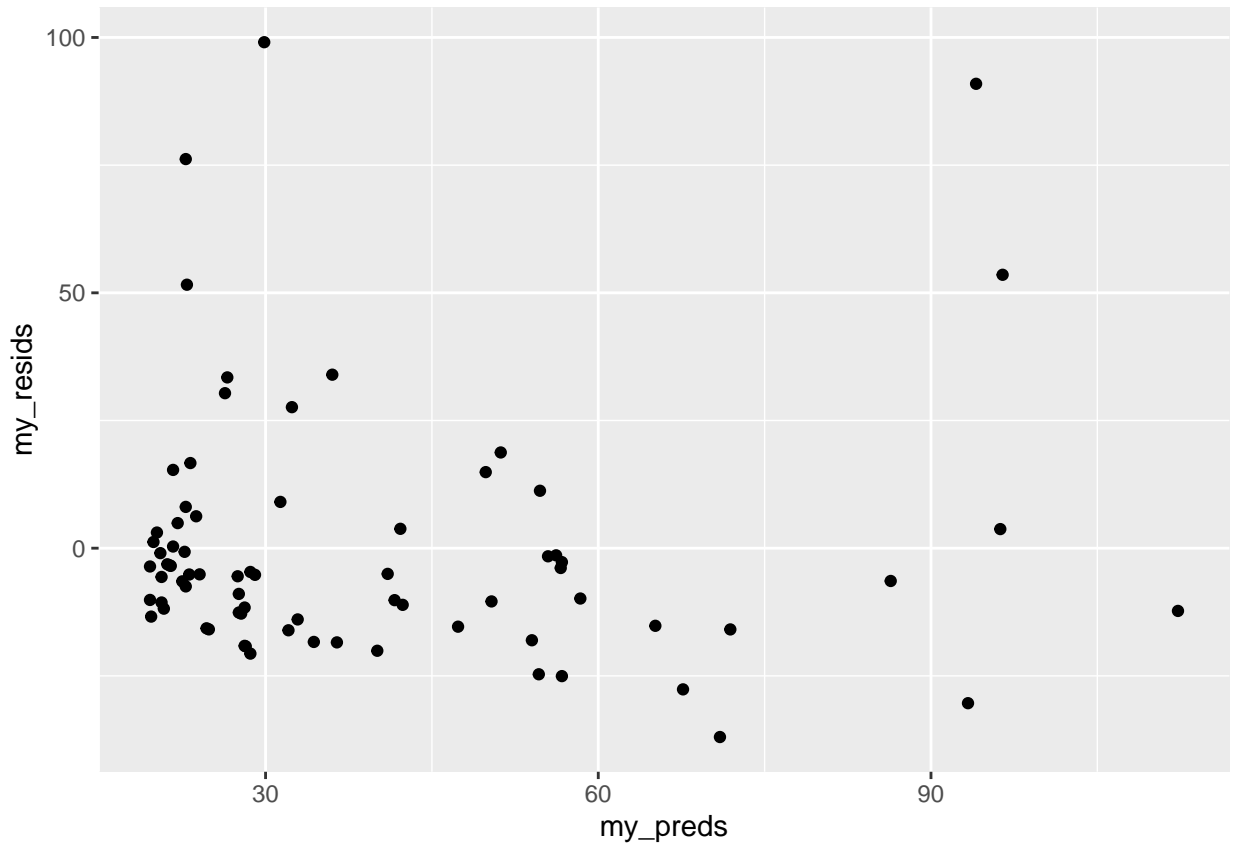
Please save your residuals and your predictions from this model as columns in the data set “legos”. The `resid()` and `predict()` functions are useful for this.

```
#DATA$NEW_VAR <- function(something)
legos$my_resids <- resid(my_line_mod)
legos$my_preds <- predict(my_line_mod)
```

Q5

Make a residual scatterplot by having the residuals of your model on the y-axis and the predicted price on the x-axis.

```
ggplot(legos,
  aes(x = my_preds,
    y = my_resids)) +
  geom_point()
```



Q6

Please comment on if the homoskedasticity and normality assumptions are met for our linear model by using the graph made in question 10.

Neither are met by a long shot. Too many outliers in the y direction, most data is negative except all the data that is past 70 which is all positive

Q7

Regardless of your answer to question 6, please write down the estimated linear regression equation. Be sure to use the name of the y and x variables in the equation and to indicate y is predicted (and not observed).

pred. amazon price = 18.96 + .104 * Pieces

Q8

Interpret your intercept from the above equation

If there are 0 pieces in the set, we predict the price to be \$19.86

Q9

Interpret your slope from the above equation

If the number of pieces in the set increase by one, we expect the mean price to increase by 10 cents

Q10

Predict the cost a lego set containing 55 pieces.

\$24.67

Q11

The Monster Truck lego set actually has 55 pieces. Using Q15 and the lego set's actual amazon price please calculate the residual. HINT: Monster Truck is the 71st row of our data set.

```
8.99 - 24.67
```

```
## [1] -15.68
```

Q12

Find R^2 . There are several ways to do this including using the `summary()` output for the model earlier or using pearson's correlation coefficient.

.4466

Q13

All said and done, do you think this model explains the relationship well?

No, the assumptions failed too hard

Transformation

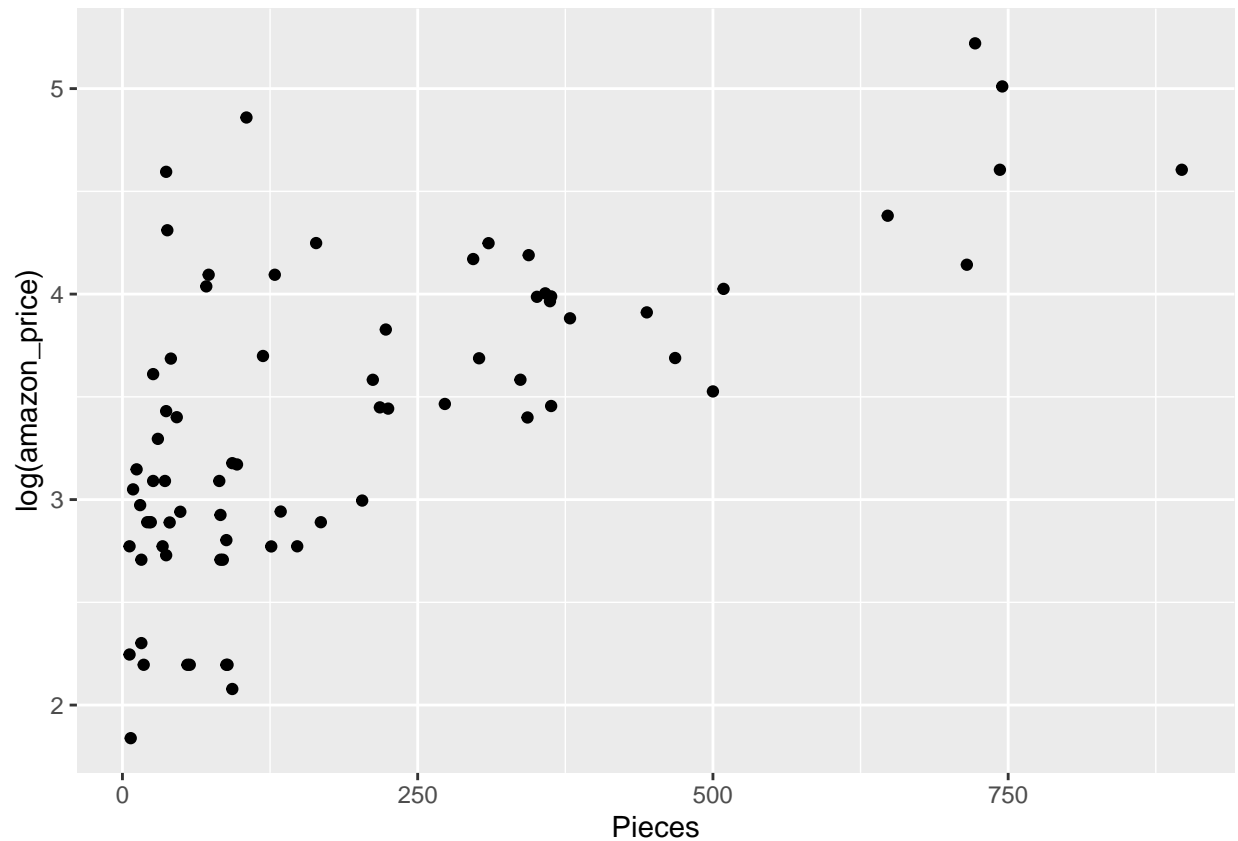
What I dislike about the residual graph I made is that there seemed to be some really outstretched values along the y-axis. That can indicate that the response variable should be transformed via a `log()` function (but not always!!).

Log-Linear Model

Q14

As such, please make a scatterplot with the log of the amazon price as the y-axis and leave the x-axis as the number of pieces used. Comment on whether you think this graph is sufficiently linear.

```
ggplot(legos,
       aes(x = Pieces,
           y = log(amazon_price))) +
  geom_point()
```

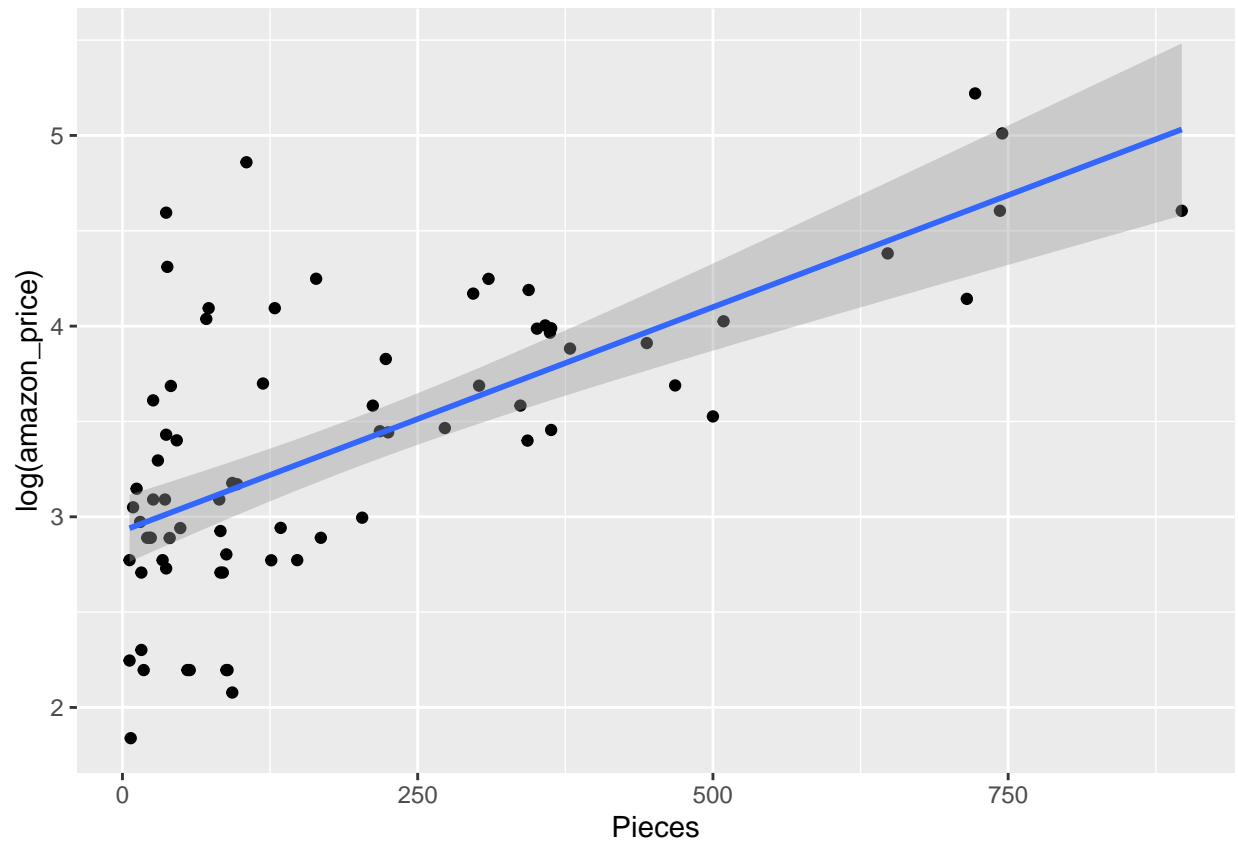


Q15

Using `log(amazon_price)` as the response variable and `pieces` as the x-axis, plot a best-fit-line using `geom_smooth`. Does the best fit line look good?

```
ggplot(legos,
       aes(x = Pieces,
           y = log(amazon_price))) +
  geom_point() +
  geom_smooth(method = 'lm')
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

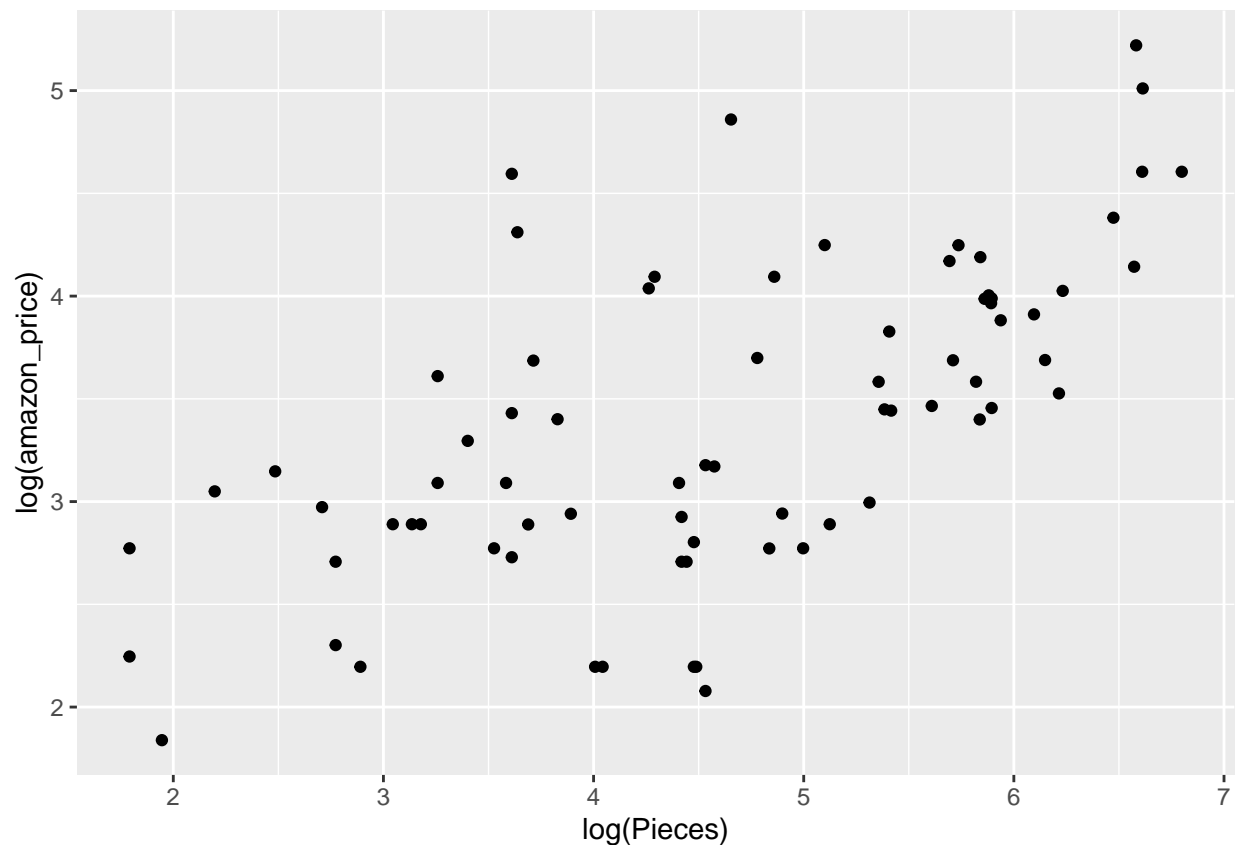
Let's try one more transformation to see if we can get something closer to what we are after

Log-Log Model

Q16

As such, please make a scatterplot with the log of the amazon price as the y-axis and the log of the number of pieces used as the x-axis. Use `geom_smooth` to fit a best-fit-line similar to question 2.

```
ggplot(legos,
  aes(x = log(Pieces),
    y = log(amazon_price))) +
  geom_point()
```



Q17

Fit a linear model using $\log(\text{amazon_price})$ as your response and $\log(\text{Pieces})$ as your explanatory variable.

```
my_log_mod <- lm(log(amazon_price) ~ log(Pieces),
                 data = legos)
```

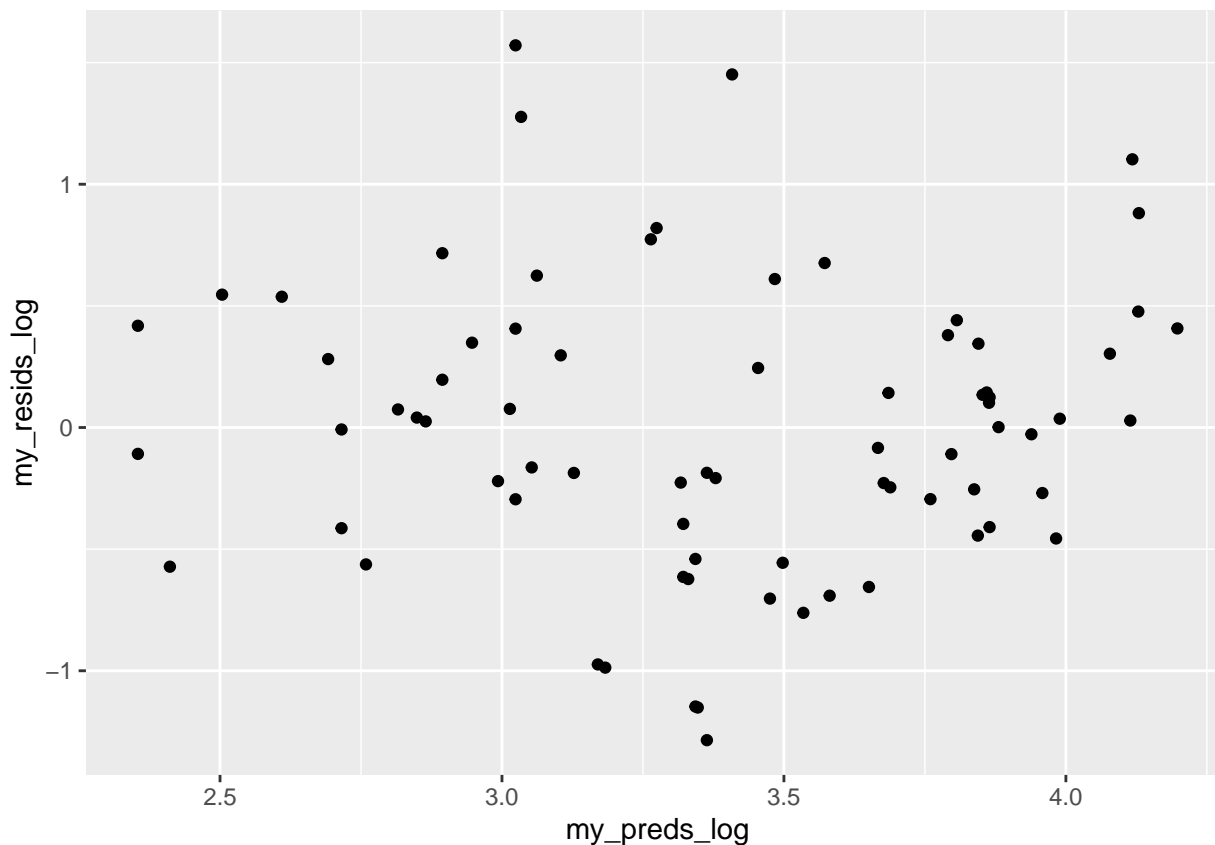
Q18

Create a residual graph for the model created in Q22 and comment on whether the normality and homoskedasticity assumptions are met.

This graph is fine, don't stare tooooo long

```
#DATA$NEW_VAR <- function(something)
legos$my_resids_log <- resid(my_log_mod)
legos$my_preds_log <- predict(my_log_mod)

ggplot(legos,
       aes(x = my_preds_log,
           y = my_resids_log)) +
  geom_point()
```



Q19

Write down your estimated equation. Be sure to indicate what the y and x variables are and that the response is estimated. Also note that in your model both variables are transformed to $\log()$'s. You do not need to back transform for this question.

predicted log amazon price = $1.69 + .368 * \log(\text{pieces})$

Q20

Interpret your value for $\hat{\beta}_0$, the intercept of your model. Be careful to differentiate between predicting the mean vs predicting the median.

If the number of pieces is set to 1 (such that the $\log(\text{pieces})$ is 0), then we expect the median price of the lego sets to be $\exp(1.69) = \$5.45$

Q21

Interpret your value for $\hat{\beta}_1$, the slope of your model. Be careful to differentiate between predicting the mean vs predicting the median.

For a 10% increase in the number of pieces we expect the median price of the lego set to increase by a multiplicative factor $1.1^{.3681} = 1.0357$ (alternatively the median increases by 3.5%)

Q22

Again, please find the predicted price for a lego set with 55 pieces using the model you just created. Be sure that the prediction is reported on the linear scale (ie I want the prediction listed in dollars). You will want to back transform for this problem.

```
log_y_hat <- 1.6948 + .3681*log(55)
exp(log_y_hat)
```

```
## [1] 23.80509
```

Q23

Using Q22's prediction, calculate the residual for the Monster Truck lego set in the data. Be sure that the residual is reported on the linear scale (ie I want the residual listed in dollars).

23.807 - 8.99 =

14.82

Indicators

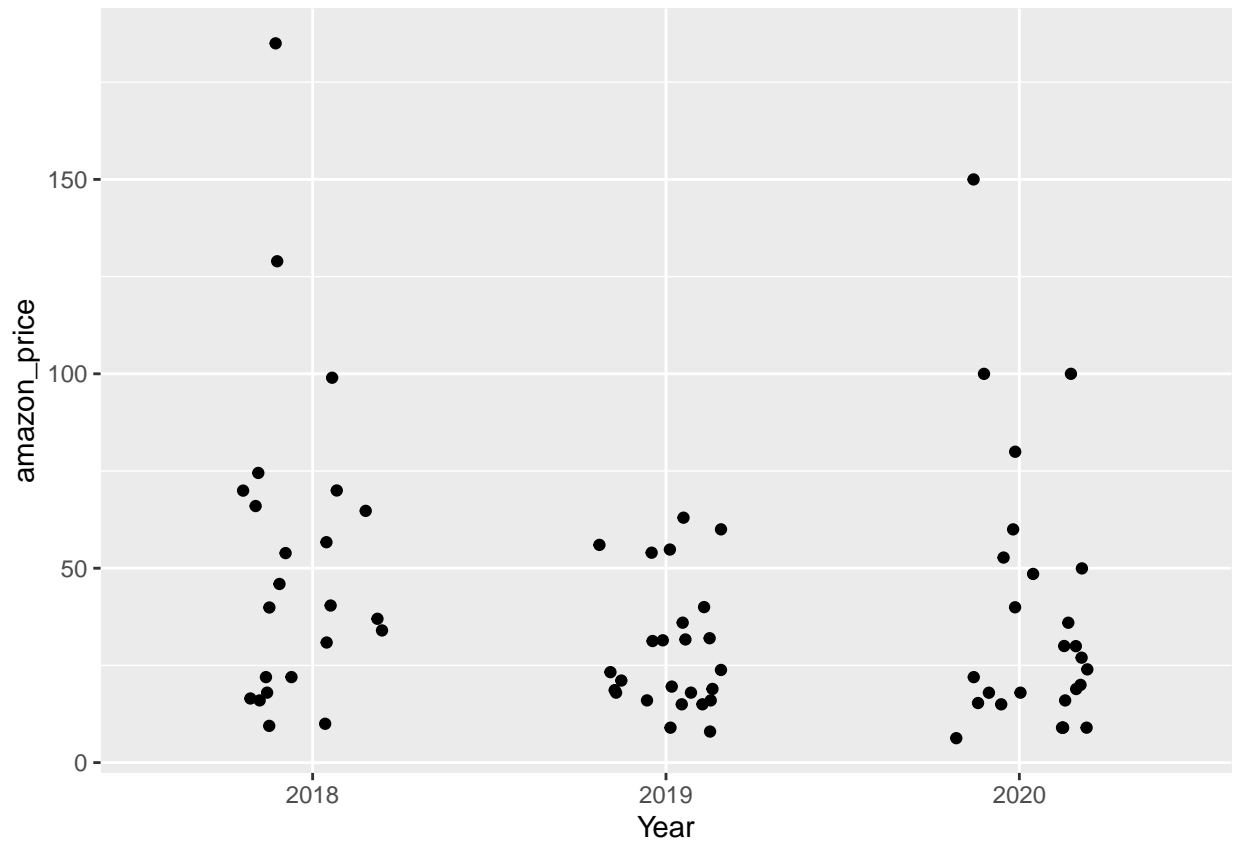
For this we are going to do something a little odd. We are going to treat Year as a categorical variable and just say that 2018, 2019, and 2020 are just labels (ie nominal) that don't mean anything numerically.

Q24

Make a plot similar to the one in the class notes. Your x-axis should be Year and your y-axis should be amazon sales price

HINT: Use `geom_jitter()` and not `geom_point()`. If the points are spread out too wide, play around with the "widths" parameter in `geom_jitter()`

```
ggplot(legos,
       aes(x = Year,
           y = amazon_price)) +
  geom_jitter(width = .2)
```



Q25

Make a linear model using Year as an explanatory variable and amazon price as the response variable.

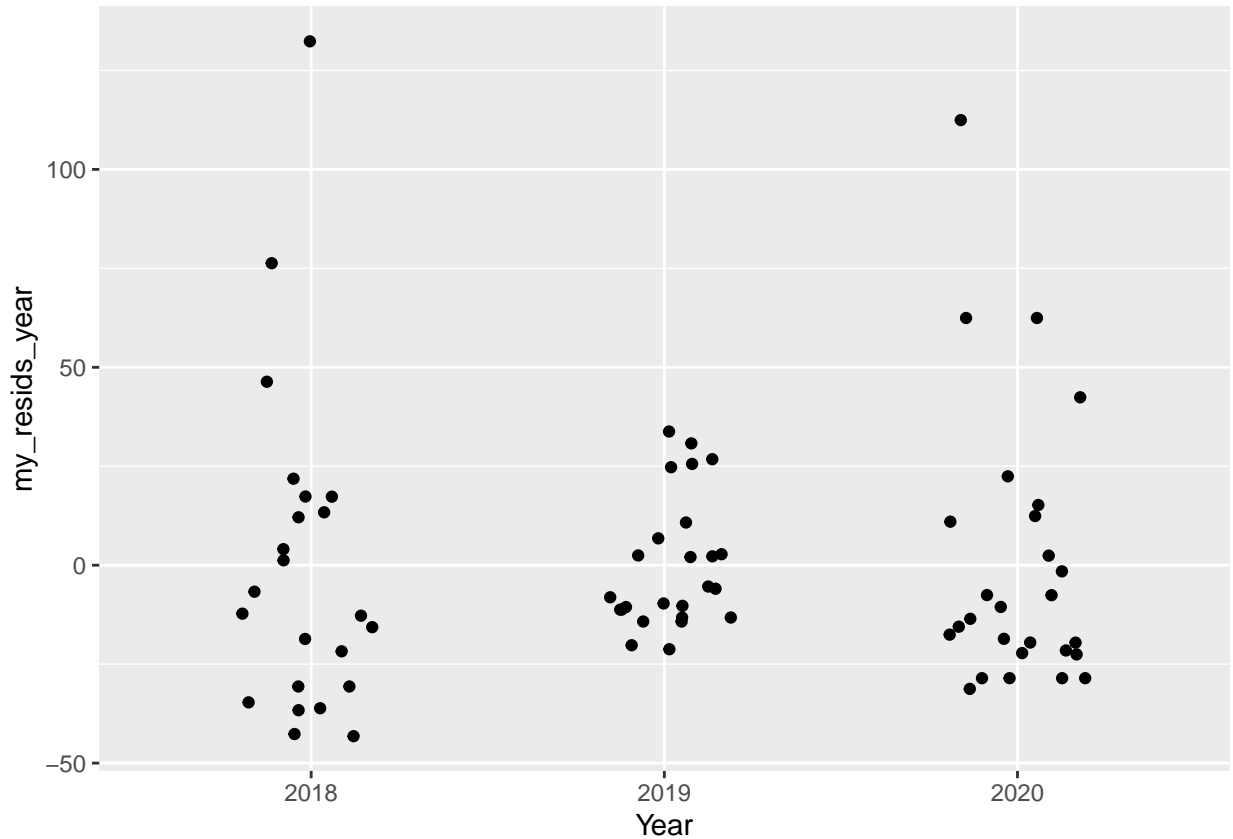
```
my_year_mod <- lm(amazon_price ~ Year,
  data = legos)
```

Q26

Make a scatterplot with your residuals on the y-axis and the x-axis being Year.

```
legos$my_resids_year <- resid(my_year_mod)
legos$my_preds_year <- predict(my_year_mod)

ggplot(legos,
  aes(x = Year,
    y = my_resids_year)) +
  geom_jitter(width = .2)
```



Q27

Comment on if the three categories (years) have heteroskedasticity or if the residuals are not normal.

HINT: Use `geom_jitter()` and not `geom_point()`. If the points are spread out too wide, play around with the “widths” parameter in `geom_jitter()`

Possibly heteroskedastic and only the year 2019 looks really normal...2020 and 2018 look like they have a right tail

Q28

Write down your best-fit-line equation. Please use the model form which uses β 's, and not the one that uses α 's. $\text{predicted amazon price} = 52.64 - 23.43 * (1_{2019}) - 15.11 * (1_{2020})$

HINT: run the `summary()` command on your model and then look at the “Coefficients” table, specifically the “Estimates” column. See the alternative slide deck for indicators for an example

Q29

Predict the cost of a lego set that was made in the 2020.

37.53

52.64 - 15.11

```
## [1] 37.53
```

Q30

Find the residual (again) for the Monster Truck data set 28.54

37.53 - 8.99

```
## [1] 28.54
```

Q31

Interpret your $\hat{\beta}_0$ value

We predict the mean price of a lego set made in 2018 to be \$52.64

Q32

Interpret your $\hat{\beta}_1$ value

We predict the change in mean price of a lego set made in 2018 to 2019 to be -\$23.42

Q33

Interpret your $\hat{\beta}_2$ value

We predict the change in mean price of a lego set made in 2018 to 2020 to be -\$15.11

Q34

Find the different between $\hat{\beta}_2$ and $\hat{\beta}_1$. What *is* this difference? What does it represent?

\$8.31, this is the change in the mean of price going from 2019 to 2020