

# homework\_2

2025-10-07

This homework will focus heavily on regression with linear-linear, log-linear, and log-log all making an appearance and indicators as well. My advice is to use the class time to focus on either the log-log models (Q18-25) or indicators (Q26 onwards). The first section is a regular linear model like you explored in the last lab.

HINT: you will make four models (one for linear, one for log-linear, one for log-log and one that deals with indicators) and I suggest you save all four models. My naming scheme is usually along the lines of

```
mod_line <- LINEAR MODEL CODE
```

```
mod_log_line <- LOG-LINEAR MODEL CODE
```

```
mod_log_log <- LOG-LOG MODEL CODE
```

```
mod_ind <- INDICATOR MODEL CODE
```

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')
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		

```
legos$Year <- as.factor(legos$Year)
```

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
- 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

## Correlation

### 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.

### Q2

Based on the three graphs in question 1, please indicate whether you think Pearson’s or Spearman’s correlation coefficients is more appropriate for talking about the 3 explanatory variable’s relationship with the response. Explain why (note: you don’t need to pick a correlation for each graph separately but just one for all three graphs)

### Q3

In your own words, please explain what it means to have a negative correlation. (No, this isn’t related to the scatterplot but it’s a good question that I want you to try to answer)

## Linear Regression

We will continue to use the lego data set for this.

### Q4

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.

### Q5

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

**Q6**

Please save your residuals and your predictions from this model. The `resid()` and `predict()` functions are useful for this.

**Q7**

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

**Q8**

Please comment on if the homoskedasticity and normality assumptions are met for our linear model by using the graph made in question 7

**Q9**

Regardless of your answer to question 8, 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).

**Q10**

Interpret your intercept from the above equation

**Q11**

Interpret your slope from the above equation

**Q12**

Predict the cost a lego set containing 55 pieces.

**Q13**

The Monster Truck lego set actually has 55 pieces. Using Q12 and the lego set's actual amazon price please calculate the residual. HINT: Monster Truck is the 71st row of our data set.

**Q14**

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. Interpret it.

**Q15**

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

## 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

#### Q16

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.

#### Q17

Using  $\log(\text{amazon\_price})$  as the response variable and pieces as the x-axis, fit a linear regression model. Plot the residuals similar to question 7 with your residuals as the y-axis and the predicted values on the x-axis. Comment if the normality and homoskedasticity assumptions are met.

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

### Log-Log Model

#### Q18

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 7.

#### Q19

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

#### Q20

Create a residual graph for the model created in Q19 and comment on whether the normality and homoskedasticity assumptions are met.

#### Q21

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.

#### Q22

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

### Q23

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

### Q24

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.

### Q25

Using Q24'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).

## 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. Making Year nominal has already been done in the code I wrote at the top of the file that reads in the data set.

### Q26

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()`

### Q27

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

### Q28

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

### Q29

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 "width" parameter in `geom_jitter()`

**Q30**

Write down your best-fit-line equation. Please use the model form which uses  $\beta$ 's, and not the one that uses  $\alpha$ 's. 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

**Q31**

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

**Q32**

Find the residual (again) for the Monster Truck set (which was made in 2020)

**Q33**

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

**Q34**

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

**Q35**

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

**Q36**

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