

READING: 01 02

EXERCISES: ALL CHAPTER 0

ASSIGNED: HW 1

PRODUCER: DR. MARIO



Course Website / Syllabus

- Access Course Website Through Canvas
- Cover Syllabus
 - Office Hours
 - Grading and Curving
 - Attendance: UNC Check-In App
 - Homework
 - Quizzes
 - Exams
 - PDFs and Gradescope
 - Grade Disputes
 - Honor Code
- Usage of Course Website and Canvas

Course Website / Syllabus

- Access Course Website Through Canvas
- Cover Syllabus
 - Office Hours
 - Grading and Curving
 - Attendance: UNC Check-In App
 - Homework
 - Quizzes
 - Exams
 - PDFs and Gradescope
 - Grade Disputes
 - Honor Code
- Usage of Course Website and Canvas



umbers with a context.

+

Departures
from a
Pattern

How do we identify the
actual pattern?

How do we characterize the
departures (errors)?

Statistical Modeling

Find a model for a relationship between a response variable (Y) and one (or more) predictor/explanatory variables (X_1, X_2, \dots, X_k).

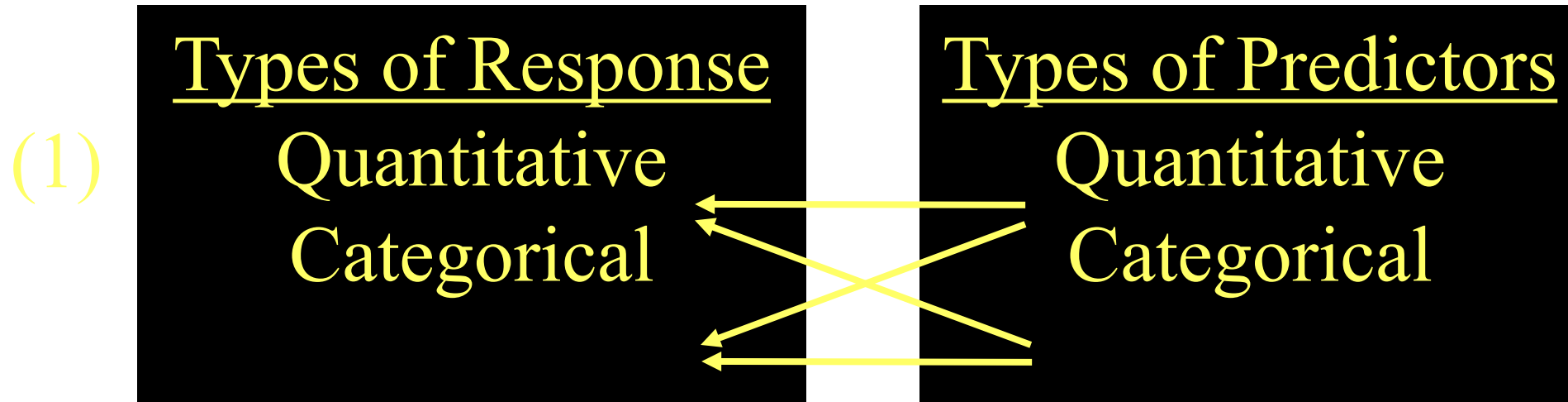
Types of Variables

Quantitative: expressible as
numbers for which arithmetic
makes sense

Categorical: divides sample points
into groups

Binary = categorical with just two groups

Two Main Themes of STOR 455



- (2) Allow for models with *multiple* predictors.

Building a Statistical Model: Four Step Process

1. CHOOSE – Pick a form for the model
2. FIT – Estimate any parameters
3. ASSESS – Is the model adequate?
Could it be simpler? Are conditions met?
4. USE – Answer the question of interest

General form of a model:

$$Y = f(X) + \varepsilon$$

Random Error

“Expected” Y for some combination of predictors

Data

Model

Error

Example: Lego Prices

Question:

How can we predict the price of a Lego set?

Predictor variables: Start with *none*.

Example: Constant Model

$$Y = c + \varepsilon \quad \text{where } c \text{ is an (unknown) constant}$$

Terminology:

The constant c is a parameter of this model.

We use data to provide a sample estimate of c .

How should we estimate c from data?

Predicted Value for Response

Get an \hat{y} for Y using the predictors and the model with estimated parameters.

Notation: The predicted y is denoted \hat{y}

For the constant model: $\hat{y} = \hat{c}$

Examples:

- \bar{y} (sample mean)
- \tilde{y} (sample median)

Residuals

Using the predicted value for each sample case the residual is

$$\text{Residual} = y - \hat{y}$$



Actual

Predicted

Technology

We need software to automate computations...

R – a free, widely used, open
source statistics package

Rstudio – an interface for R

RStudio

The screenshot shows the RStudio application window. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. Below the menu is a toolbar with icons for file operations and a search bar. The main editor area on the left contains a code chunk with R Markdown syntax: `---`, `title: "R Notebook"`, `output: html_notebook`, `---`, followed by a comment and a code chunk `{r}` containing `plot(cars)`. A yellow annotation box over the editor says "Editor: write/view code, data".

On the right side, the Environment pane shows the Global Environment with a yellow annotation box stating "Environment: lists active variables, functions, ...". Below it, the Files pane shows the project directory structure with a yellow annotation box listing: "Access files", "View plots", "Control packages", and "View help".

At the bottom, the Console pane shows the R prompt and some initial output, with a yellow annotation box stating "Console: Enter commands, view output, error messages".

A First RStudio Session

- Load the *Lego* data into R
- Summarize the *Amazon_Price* variable
 - Numerical: mean and median
 - Graphical: histogram, boxplot
- Compute and evaluate residuals

A First RStudio Session

Load the *lego* data into R

```
` ``{r}
# loads a package needed to use the read_csv() function
# install package before first using it for the first time

library(readr)

# loads the lego dataframe into the environment from GitHub

lego <- read_csv("https://raw.githubusercontent.com/JA-McLean/STOR455/master/data/lego.csv")

# Alternative way to load dataframe (remove # to use)
# lego.csv must be saved in the same folder as this notebook!

#lego <- read_csv("lego.csv")

# Shows the variables and first 6 cases (by default)

head(lego)
` ``
```

A First RStudio Session

Summarize the *Amazon_Price* variable - Numerical: mean and median

```
# dataframe$variable_name
```

```
mean(lego$Amazon_Price, na.rm = TRUE)
```

```
median(lego$Amazon_Price, na.rm = TRUE)
```

```
[1] 57.8232
```

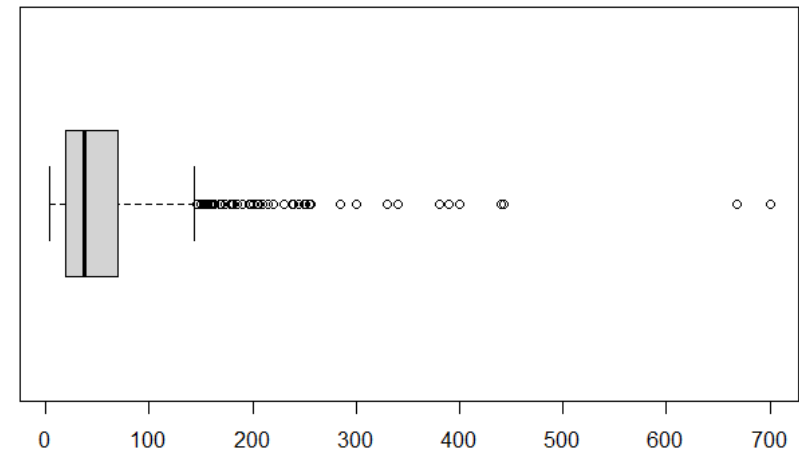
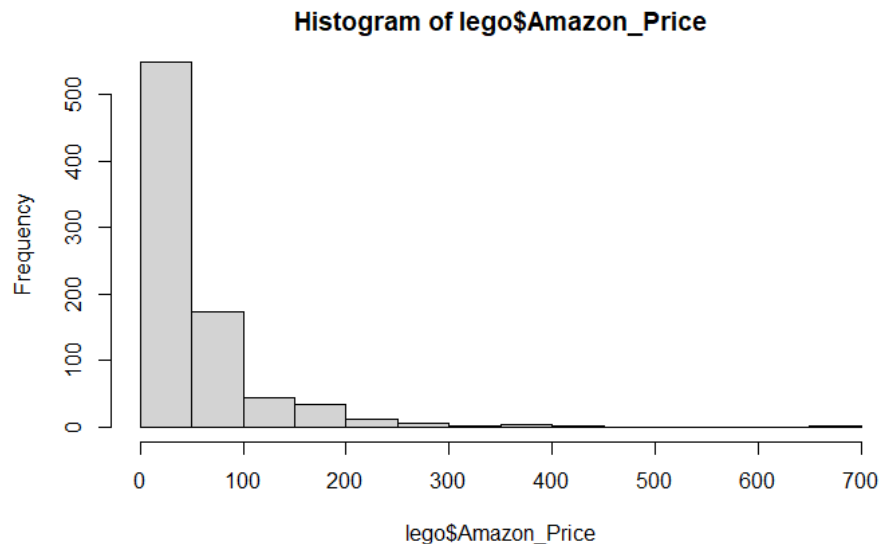
```
[1] 37.325
```

A First RStudio Session

Summarize the *Amazon_Price* variable - Graphical: histogram, boxplot

```
hist(lego$Amazon_Price)
```

```
boxplot(lego$Amazon_Price, horizontal = TRUE)
```



A First RStudio Session

Compute and evaluate residuals

```
# removes NA Amazon_Prices

lego_rm_AP_na = subset(lego, is.na(Amazon_Price) == FALSE)

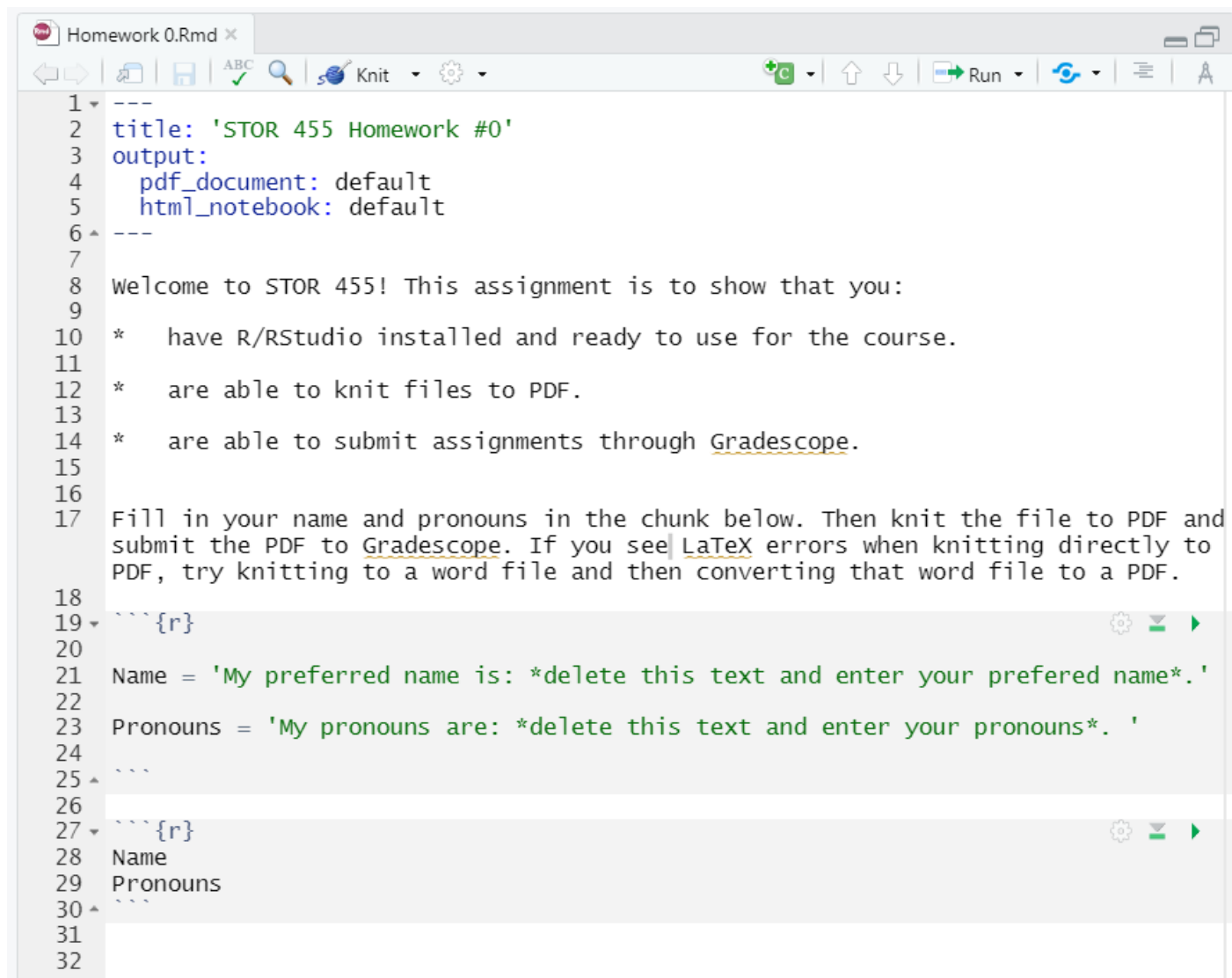
# Assignment operators in R: = vs. <-

xbar = mean(lego_rm_AP_na$Amazon_Price)
m = median(lego_rm_AP_na$Amazon_Price)

residxbar = lego_rm_AP_na$Amazon_Price - xbar
residm = lego_rm_AP_na$Amazon_Price - m

sum(residxbar^2)
sum(residm^2)
```


Homework 0



```
1 ---
2 title: 'STOR 455 Homework #0'
3 output:
4   pdf_document: default
5   html_notebook: default
6 ---
7
8 Welcome to STOR 455! This assignment is to show that you:
9
10 * have R/RStudio installed and ready to use for the course.
11
12 * are able to knit files to PDF.
13
14 * are able to submit assignments through Gradescope.
15
16
17 Fill in your name and pronouns in the chunk below. Then knit the file to PDF and
18 submit the PDF to Gradescope. If you see LaTeX errors when knitting directly to
19 PDF, try knitting to a word file and then converting that word file to a PDF.
20
21 ```{r}
22
23 Name = 'My preferred name is: *delete this text and enter your preferred name*.'
24
25 Pronouns = 'My pronouns are: *delete this text and enter your pronouns*.'
26
27 ```
28
29 Name
30 Pronouns
31
32 ```
```